

Derivations of Applied Mathematics

Thaddeus H. Black

Revised 10 November 2007

Thaddeus H. Black, 1967–.
Derivations of Applied Mathematics.
10 November 2007.
U.S. Library of Congress class QA401.

Copyright © 1983–2007 by Thaddeus H. Black <thb@derivations.org>.

Published by the Debian Project [8].

This book is free software. You can redistribute and/or modify it under the terms of the GNU General Public License [14], version 2.

Contents

Preface	xv
1 Introduction	1
1.1 Applied mathematics	1
1.2 Rigor	2
1.2.1 Axiom and definition	2
1.2.2 Mathematical extension	4
1.3 Complex numbers and complex variables	5
1.4 On the text	5
2 Classical algebra and geometry	7
2.1 Basic arithmetic relationships	7
2.1.1 Commutivity, associativity, distributivity	7
2.1.2 Negative numbers	9
2.1.3 Inequality	10
2.1.4 The change of variable	10
2.2 Quadratics	11
2.3 Notation for series sums and products	13
2.4 The arithmetic series	15
2.5 Powers and roots	15
2.5.1 Notation and integral powers	16
2.5.2 Roots	17
2.5.3 Powers of products and powers of powers	19
2.5.4 Sums of powers	20
2.5.5 Summary and remarks	20
2.6 Multiplying and dividing power series	21
2.6.1 Multiplying power series	22
2.6.2 Dividing power series	22
2.6.3 Dividing power series by matching coefficients	26

2.6.4	Common quotients and the geometric series	29
2.6.5	Variations on the geometric series	29
2.7	Constants and variables	30
2.8	Exponentials and logarithms	32
2.8.1	The logarithm	32
2.8.2	Properties of the logarithm	33
2.9	Triangles and other polygons: simple facts	33
2.9.1	Triangle area	34
2.9.2	The triangle inequalities	34
2.9.3	The sum of interior angles	35
2.10	The Pythagorean theorem	36
2.11	Functions	38
2.12	Complex numbers (introduction)	39
2.12.1	Rectangular complex multiplication	41
2.12.2	Complex conjugation	41
2.12.3	Power series and analytic functions (preview)	43
3	Trigonometry	45
3.1	Definitions	45
3.2	Simple properties	47
3.3	Scalars, vectors, and vector notation	47
3.4	Rotation	51
3.5	Trigonometric sums and differences	53
3.5.1	Variations on the sums and differences	54
3.5.2	Trigonometric functions of double and half angles	55
3.6	Trigonometrics of the hour angles	55
3.7	The laws of sines and cosines	59
3.8	Summary of properties	60
3.9	Cylindrical and spherical coordinates	62
3.10	The complex triangle inequalities	64
3.11	De Moivre's theorem	65
4	The derivative	67
4.1	Infinitesimals and limits	67
4.1.1	The infinitesimal	68
4.1.2	Limits	69
4.2	Combinatorics	70
4.2.1	Combinations and permutations	70
4.2.2	Pascal's triangle	72
4.3	The binomial theorem	72

4.3.1	Expanding the binomial	72
4.3.2	Powers of numbers near unity	73
4.3.3	Complex powers of numbers near unity	74
4.4	The derivative	75
4.4.1	The derivative of the power series	75
4.4.2	The Leibnitz notation	76
4.4.3	The derivative of a function of a complex variable	78
4.4.4	The derivative of z^a	79
4.4.5	The logarithmic derivative	80
4.5	Basic manipulation of the derivative	80
4.5.1	The derivative chain rule	81
4.5.2	The derivative product rule	81
4.6	Extrema and higher derivatives	83
4.7	L'Hôpital's rule	84
4.8	The Newton-Raphson iteration	86
5	The complex exponential	91
5.1	The real exponential	91
5.2	The natural logarithm	94
5.3	Fast and slow functions	96
5.4	Euler's formula	98
5.5	Complex exponentials and de Moivre	101
5.6	Complex trigonometrics	101
5.7	Summary of properties	103
5.8	Derivatives of complex exponentials	103
5.8.1	Derivatives of sine and cosine	103
5.8.2	Derivatives of the trigonometrics	106
5.8.3	Derivatives of the inverse trigonometrics	107
5.9	The actuality of complex quantities	108
6	Primes, roots and averages	111
6.1	Prime numbers	111
6.1.1	The infinite supply of primes	111
6.1.2	Compositional uniqueness	112
6.1.3	Rational and irrational numbers	115
6.2	The existence and number of roots	116
6.2.1	Polynomial roots	116
6.2.2	The fundamental theorem of algebra	117
6.3	Addition and averages	118
6.3.1	Serial and parallel addition	118

6.3.2	Averages	121
7	The integral	125
7.1	The concept of the integral	125
7.1.1	An introductory example	126
7.1.2	Generalizing the introductory example	129
7.1.3	The balanced definition and the trapezoid rule	129
7.2	The antiderivative	130
7.3	Operators, linearity and multiple integrals	132
7.3.1	Operators	132
7.3.2	A formalism	133
7.3.3	Linearity	134
7.3.4	Summational and integrodifferential transitivity	135
7.3.5	Multiple integrals	137
7.4	Areas and volumes	139
7.4.1	The area of a circle	139
7.4.2	The volume of a cone	139
7.4.3	The surface area and volume of a sphere	140
7.5	Checking an integration	143
7.6	Contour integration	144
7.7	Discontinuities	145
7.8	Remarks (and exercises)	148
8	The Taylor series	151
8.1	The power-series expansion of $1/(1 - z)^{n+1}$	152
8.1.1	The formula	152
8.1.2	The proof by induction	154
8.1.3	Convergence	155
8.1.4	General remarks on mathematical induction	156
8.2	Shifting a power series' expansion point	157
8.3	Expanding functions in Taylor series	159
8.4	Analytic continuation	161
8.5	Branch points	163
8.6	Entire and meromorphic functions	165
8.7	Extrema over a complex domain	166
8.8	Cauchy's integral formula	167
8.8.1	The meaning of the symbol dz	168
8.8.2	Integrating along the contour	168
8.8.3	The formula	172
8.8.4	Enclosing a multiple pole	173

8.9	Taylor series for specific functions	174
8.10	Error bounds	175
8.10.1	Examples	175
8.10.2	Majorization	178
8.10.3	Geometric majorization	179
8.10.4	Calculation outside the fast convergence domain . . .	181
8.10.5	Nonconvergent series	183
8.10.6	Remarks	184
8.11	Calculating 2π	185
8.12	Odd and even functions	185
8.13	Trigonometric poles	186
8.14	The Laurent series	187
8.15	Taylor series in $1/z$	190
8.16	The multidimensional Taylor series	191
9	Integration techniques	193
9.1	Integration by antiderivative	193
9.2	Integration by substitution	194
9.3	Integration by parts	195
9.4	Integration by unknown coefficients	197
9.5	Integration by closed contour	200
9.6	Integration by partial-fraction expansion	205
9.6.1	Partial-fraction expansion	205
9.6.2	Repeated poles	206
9.6.3	Integrating a rational function	209
9.6.4	The derivatives of a rational function	210
9.6.5	Repeated poles (the conventional technique)	210
9.6.6	The existence and uniqueness of solutions	212
9.7	Frullani's integral	213
9.8	Products of exponentials, powers and logs	215
9.9	Integration by Taylor series	216
10	Cubics and quartics	219
10.1	Vieta's transform	220
10.2	Cubics	220
10.3	Superfluous roots	223
10.4	Edge cases	225
10.5	Quartics	227
10.6	Guessing the roots	229

11 The matrix	233
11.1 Provenance and basic use	236
11.1.1 The linear transformation	236
11.1.2 Matrix multiplication (and addition)	237
11.1.3 Row and column operators	239
11.1.4 The transpose and the adjoint	240
11.2 The Kronecker delta	241
11.3 Dimensionality and matrix forms	242
11.3.1 The null and dimension-limited matrices	244
11.3.2 The identity, scalar and extended matrices	245
11.3.3 The active region	247
11.3.4 Other matrix forms	247
11.3.5 The rank- r identity matrix	248
11.3.6 The truncation operator	249
11.3.7 The elementary vector and lone-element matrix	249
11.3.8 Off-diagonal entries	250
11.4 The elementary operator	250
11.4.1 Properties	252
11.4.2 Commutation and sorting	253
11.5 Inversion and similarity (introduction)	254
11.6 Parity	258
11.7 The quasidelementary operator	260
11.7.1 The general interchange operator	261
11.7.2 The general scaling operator	262
11.7.3 Addition quasidelementaries	263
11.8 The unit triangular matrix	265
11.8.1 Construction	267
11.8.2 The product of like unit triangular matrices	268
11.8.3 Inversion	268
11.8.4 The parallel unit triangular matrix	269
11.8.5 The partial unit triangular matrix	274
11.9 The shift operator	276
11.10 The Jacobian derivative	276
12 Rank and the Gauss-Jordan	279
12.1 Linear independence	280
12.2 The elementary similarity transformation	282
12.3 The Gauss-Jordan decomposition	282
12.3.1 Motive	284
12.3.2 Method	287

12.3.3	The algorithm	288
12.3.4	Rank and independent rows	296
12.3.5	Inverting the factors	297
12.3.6	Truncating the factors	298
12.3.7	Properties of the factors	299
12.3.8	Marginalizing the factor I_n	300
12.3.9	Decomposing an extended operator	301
12.4	Vector replacement	302
12.5	Rank	305
12.5.1	A logical maneuver	306
12.5.2	The impossibility of identity-matrix promotion	307
12.5.3	General matrix rank and its uniqueness	310
12.5.4	The full-rank matrix	312
12.5.5	Under- and overdetermined systems (introduction)	313
12.5.6	The full-rank factorization	314
12.5.7	Full column rank and the Gauss-Jordan's K and S	315
12.5.8	The significance of rank uniqueness	316
13	Inversion and orthonormalization	319
13.1	Inverting the square matrix	320
13.2	The exactly determined linear system	324
13.3	The kernel	325
13.3.1	The Gauss-Jordan kernel formula	326
13.3.2	Converting between kernel matrices	330
13.3.3	The degree of freedom	330
13.4	The nonoverdetermined linear system	332
13.4.1	Particular and homogeneous solutions	333
13.4.2	A particular solution	333
13.4.3	The general solution	334
13.5	The residual	334
13.6	The pseudoinverse and least squares	335
13.6.1	Least squares in the real domain	337
13.6.2	The invertibility of A^*A	338
13.6.3	Positive definiteness	339
13.6.4	The Moore-Penrose pseudoinverse	339
13.7	The multivariate Newton-Raphson iteration	343
13.8	The dot product	344
13.9	The orthogonal complement	346
13.10	Gram-Schmidt orthonormalization	346
13.10.1	Efficient implementation	348

13.10.2 The Gram-Schmidt decomposition	348
13.10.3 The Gram-Schmidt kernel formula	352
13.11 The unitary matrix	354
14 The eigenvalue	357
14.1 The determinant	357
14.1.1 Basic properties	359
14.1.2 The determinant and the elementary operator	362
14.1.3 The determinant of a singular matrix	363
14.1.4 The determinant of a matrix product	364
14.1.5 Determinants of inverse and unitary matrices	364
14.1.6 Inverting the square matrix by determinant	365
14.2 Coincident properties	366
14.3 The eigenvalue itself	368
14.4 The eigenvector	369
14.5 Eigensolution facts	370
14.6 Diagonalization	371
14.7 Remarks on the eigenvalue	373
14.8 Matrix condition	374
14.9 The similarity transformation	375
14.10 The Schur decomposition	377
14.10.1 Derivation	377
14.10.2 The nondiagonalizable matrix	382
14.11 The Hermitian matrix	384
14.12 The singular value decomposition	388
14.13 General remarks on the matrix	390
15 Vector algebra (to be written)	393
A Hex and other notational matters	397
A.1 Hexadecimal numerals	398
A.2 Avoiding notational clutter	399
B The Greek alphabet	401
C Manuscript history	405

List of Tables

2.1	Basic properties of arithmetic.	8
2.2	Power properties and definitions.	16
2.3	Dividing power series through successively smaller powers. . .	25
2.4	Dividing power series through successively larger powers. . .	27
2.5	General properties of the logarithm.	34
3.1	Simple properties of the trigonometric functions.	48
3.2	Trigonometric functions of the hour angles.	58
3.3	Further properties of the trigonometric functions.	61
3.4	Rectangular, cylindrical and spherical coordinate relations. .	63
5.1	Complex exponential properties.	104
5.2	Derivatives of the trigonometrics.	106
5.3	Derivatives of the inverse trigonometrics.	108
6.1	Parallel and serial addition identities.	120
7.1	Basic derivatives for the antiderivative.	132
8.1	Taylor series.	176
9.1	Antiderivatives of products of exps, powers and logs. . . .	216
10.1	The method to extract the three roots of the general cubic. .	223
10.2	The method to extract the four roots of the general quartic. .	230
11.1	Elementary operators: interchange.	255
11.2	Elementary operators: scaling.	256
11.3	Elementary operators: addition.	256
11.4	Matrix inversion properties.	258
11.5	Properties of the parallel unit triangular matrix.	273

12.1	Some elementary similarity transformations.	283
12.2	A few properties of the Gauss-Jordan factors.	300
12.3	The symmetrical equations of § 12.4.	304
B.1	The Roman and Greek alphabets.	402

List of Figures

1.1	Two triangles.	4
2.1	Multiplicative commutivity.	8
2.2	The sum of a triangle's inner angles: turning at the corner. .	35
2.3	A right triangle.	37
2.4	The Pythagorean theorem.	37
2.5	The complex (or Argand) plane.	40
3.1	The sine and the cosine.	46
3.2	The sine function.	47
3.3	A two-dimensional vector $\mathbf{u} = \hat{\mathbf{x}}x + \hat{\mathbf{y}}y$	49
3.4	A three-dimensional vector $\mathbf{v} = \hat{\mathbf{x}}x + \hat{\mathbf{y}}y + \hat{\mathbf{z}}z$	49
3.5	Vector basis rotation.	52
3.6	The 0x18 hours in a circle.	57
3.7	Calculating the hour trigonometrics.	57
3.8	The laws of sines and cosines.	59
3.9	A point on a sphere.	63
4.1	The plan for Pascal's triangle.	72
4.2	Pascal's triangle.	73
4.3	A local extremum.	83
4.4	A level inflection.	84
4.5	The Newton-Raphson iteration.	87
5.1	The natural exponential.	94
5.2	The natural logarithm.	95
5.3	The complex exponential and Euler's formula.	99
5.4	The derivatives of the sine and cosine functions.	105
7.1	Areas representing discrete sums.	126

7.2	An area representing an infinite sum of infinitesimals.	128
7.3	Integration by the trapezoid rule.	130
7.4	The area of a circle.	139
7.5	The volume of a cone.	140
7.6	A sphere.	141
7.7	An element of a sphere's surface.	141
7.8	A contour of integration.	145
7.9	The Heaviside unit step $u(t)$	146
7.10	The Dirac delta $\delta(t)$	146
8.1	A complex contour of integration in two segments.	169
8.2	A Cauchy contour integral.	173
8.3	Majorization.	179
9.1	Integration by closed contour.	201
10.1	Vieta's transform, plotted logarithmically.	221
13.1	Fitting a line to measured data.	336

Preface

I never meant to write this book. It emerged unheralded, unexpectedly.

The book began in 1983 when a high-school classmate challenged me to prove the Pythagorean theorem on the spot. I lost the dare, but looking the proof up later I recorded it on loose leaves, adding to it the derivations of a few other theorems of interest to me. From such a kernel the notes grew over time, until family and friends suggested that the notes might make the material for a book.

The book, a work yet in progress but a complete, entire book as it stands, first frames coherently the simplest, most basic derivations of applied mathematics, treating quadratics, trigonometrics, exponentials, derivatives, integrals, series, complex variables and, of course, the aforementioned Pythagorean theorem. These and others establish the book's foundation in Chapters 2 through 9. Later chapters build upon the foundation, deriving results less general or more advanced. Such is the book's plan.

The book is neither a tutorial on the one hand nor a bald reference on the other. It is a study reference, in the tradition of, for instance, Kernighan's and Ritchie's *The C Programming Language* [23]. In the book, you can look up some particular result directly, or you can begin on page one and read—with toil and commensurate profit—straight through to the end of the last chapter. The book is so organized.

The book generally follows established applied mathematical convention but where worthwhile occasionally departs therefrom. One particular respect in which the book departs requires some defense here, I think: the book employs hexadecimal numerals.

There is nothing wrong with decimal numerals as such. Decimal numerals are fine in history and anthropology (man has ten fingers), finance and accounting (dollars, cents, pounds, shillings, pence: the base hardly matters), law and engineering (the physical units are arbitrary anyway); but they are merely serviceable in mathematical theory, never aesthetic. There unfortunately really is no gradual way to bridge the gap to hexadecimal

(shifting to base eleven, thence to twelve, etc., is no use). If one wishes to reach hexadecimal ground, one must leap. Twenty years of keeping my own private notes in hex have persuaded me that the leap justifies the risk. In other matters, by contrast, the book leaps seldom. The book in general walks a tolerably conventional applied mathematical line.

The book belongs to the emerging tradition of open-source software, where at the time of this writing it fills a void. Nevertheless it is a *book*, not a program. Lore among open-source developers holds that open development inherently leads to superior work. Well, maybe. Often it does in fact. Personally with regard to my own work, I should rather not make too many claims. It would be vain to deny that professional editing and formal peer review, neither of which the book enjoys, had substantial value. On the other hand, it does not do to despise the *amateur* (literally, one who does for the love of it: not such a bad motive, after all¹) on principle, either—unless one would on the same principle despise a Washington or an Einstein, or a Debian Developer [8]. Open source has a spirit to it which leads readers to be far more generous with their feedback than ever could be the case with a traditional, proprietary book. Such readers, among whom a surprising concentration of talent and expertise are found, enrich the work freely. This has value, too.

The book's peculiar mission and program lend it an unusual quantity of discursive footnotes. These footnotes offer nonessential material which, while edifying, coheres insufficiently well to join the main narrative. The footnote is an imperfect messenger, of course. Catching the reader's eye, it can break the flow of otherwise good prose. Modern publishing offers various alternatives to the footnote—numbered examples, sidebars, special fonts, colored inks, etc. Some of these are merely trendy. Others, like numbered examples, really do help the right kind of book; but for this book the humble footnote, long sanctioned by an earlier era of publishing, extensively employed by such sages as Gibbon [16] and Shirer [38], seems the most able messenger. In this book it shall have many messages to bear.

The book naturally subjoins an index. The canny reader will avoid using the index (of this and most other books), tending rather to consult the table of contents. In typical science/engineering style, the book numbers its sections, figures, tables and formulas, but not its theorems, the last of which it generally sets in italic type. A theorem is referenced by the number of the section that states it.

The book provides a bibliography listing other books I have referred to

¹The expression is derived from an observation I seem to recall George F. Will making.

while writing. Mathematics by its very nature promotes queer bibliographies, however, for its methods and truths are established by derivation rather than authority. Much of the book consists of common mathematical knowledge or of proofs I have worked out with my own pencil from various ideas gleaned who knows where over the years. The latter proofs are perhaps original or semi-original from my personal point of view, but it is unlikely that many if any of them are truly new. To the initiated, the mathematics itself often tends to suggest the form of the proof: if to me, then surely also to others who came before; and even where a proof is new the idea proven probably is not.

Among the bibliography's entries stands a reference [5] to my doctoral adviser G.S. Brown, though the book's narrative seldom explicitly invokes the reference. Prof. Brown had nothing directly to do with the book's development, for a substantial bulk of the manuscript, or of the notes underlying it, had been drafted years before I had ever heard of Prof. Brown or he of me, and my work under him did not regard the book in any case. However, the ways in which a good doctoral adviser influences his student are both complex and deep. Prof. Brown's style and insight touch the book in many places and in many ways, usually too implicitly coherently to cite.

As to a grand goal, underlying purpose or hidden plan, the book has none, other than to derive as many useful mathematical results as possible and to record the derivations together in an orderly manner in a single volume. What constitutes "useful" or "orderly" is a matter of perspective and judgment, of course. My own peculiar heterogeneous background in military service, building construction, electrical engineering, electromagnetic analysis and Debian development, my nativity, residence and citizenship in the United States, undoubtedly bias the selection and presentation to some degree. How other authors go about writing their books, I do not know, but I suppose that what is true for me is true for many of them also: we begin by organizing notes for our own use, then observe that the same notes might prove useful to others, and then undertake to revise the notes and to bring them into a form which actually is useful to others. Whether this book succeeds in the last point is for the reader to judge.

THB

Chapter 1

Introduction

This is a book of applied mathematical proofs. If you have seen a mathematical result, if you want to know why the result is so, you can look for the proof here.

The book's purpose is to convey the essential ideas underlying the derivations of a large number of mathematical results useful in the modeling of physical systems. To this end, the book emphasizes main threads of mathematical argument and the motivation underlying the main threads, deëmphasizing formal mathematical rigor. It derives mathematical results from the purely applied perspective of the scientist and the engineer.

The book's chapters are topical. This first chapter treats a few introductory matters of general interest.

1.1 Applied mathematics

What is applied mathematics?

Applied mathematics is a branch of mathematics that concerns itself with the application of mathematical knowledge to other domains.... The question of what is applied mathematics does not answer to logical classification so much as to the sociology of professionals who use mathematics. [26]

That is about right, on both counts. In this book we shall define *applied mathematics* to be correct mathematics useful to scientists, engineers and the like; proceeding not from reduced, well defined sets of axioms but rather directly from a nebulous mass of natural arithmetical, geometrical and classical-algebraic idealizations of physical systems; demonstrable but generally lacking the detailed rigor of the professional mathematician.

1.2 Rigor

It is impossible to write such a book as this without some discussion of mathematical rigor. Applied and professional mathematics differ principally and essentially in the layer of abstract definitions the latter subimposes beneath the physical ideas the former seeks to model. Notions of mathematical rigor fit far more comfortably in the abstract realm of the professional mathematician; they do not always translate so gracefully to the applied realm. The applied mathematical reader and practitioner needs to be aware of this difference.

1.2.1 Axiom and definition

Ideally, a professional mathematician knows or precisely specifies in advance the set of fundamental axioms he means to use to derive a result. A prime aesthetic here is irreducibility: no axiom in the set should overlap the others or be specifiable in terms of the others. Geometrical argument—proof by sketch—is distrusted. The professional mathematical literature discourages undue pedantry indeed, but its readers do implicitly demand a convincing assurance that its writers *could* derive results in pedantic detail if called upon to do so. Precise definition here is critically important, which is why the professional mathematician tends not to accept blithe statements such as that

$$\frac{1}{0} = \infty,$$

without first inquiring as to exactly what is meant by symbols like 0 and ∞ .

The applied mathematician begins from a different base. His ideal lies not in precise definition or irreducible axiom, but rather in the elegant modeling of the essential features of some physical system. Here, mathematical definitions tend to be made up *ad hoc* along the way, based on previous experience solving similar problems, adapted implicitly to suit the model at hand. If you ask the applied mathematician exactly what his axioms are, which symbolic algebra he is using, he usually doesn't know; what he knows is that the bridge has its footings in certain soils with specified tolerances, suffers such-and-such a wind load, etc. To avoid error, the applied mathematician relies not on abstract formalism but rather on a thorough mental grasp of the essential physical features of the phenomenon he is trying to model. An equation like

$$\frac{1}{0} = \infty$$

may make perfect sense without further explanation to an applied mathematical readership, depending on the physical context in which the equation is introduced. Geometrical argument—proof by sketch—is not only trusted but treasured. Abstract definitions are wanted only insofar as they smooth the analysis of the particular physical problem at hand; such definitions are seldom promoted for their own sakes.

The irascible Oliver Heaviside, responsible for the important applied mathematical technique of phasor analysis, once said,

It is shocking that young people should be addling their brains over mere logical subtleties, trying to understand the proof of one obvious fact in terms of something equally ... obvious. [31]

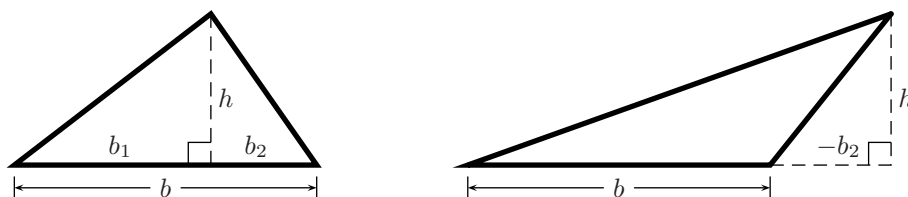
Exaggeration, perhaps, but from the applied mathematical perspective Heaviside nevertheless had a point. The professional mathematicians Richard Courant and David Hilbert put it more soberly in 1924 when they wrote,

Since the seventeenth century, physical intuition has served as a vital source for mathematical problems and methods. Recent trends and fashions have, however, weakened the connection between mathematics and physics; mathematicians, turning away from the roots of mathematics in intuition, have concentrated on refinement and emphasized the postulational side of mathematics, and at times have overlooked the unity of their science with physics and other fields. In many cases, physicists have ceased to appreciate the attitudes of mathematicians. [6, Preface]

Although the present book treats “the attitudes of mathematicians” with greater deference than some of the unnamed 1924 physicists might have done, still, Courant and Hilbert could have been speaking for the engineers and other applied mathematicians of our own day as well as for the physicists of theirs. To the applied mathematician, the mathematics is not principally meant to be developed and appreciated for its own sake; it is meant to be *used*. This book adopts the Courant-Hilbert perspective.

The introduction you are now reading is not the right venue for an essay on why both kinds of mathematics—applied and professional (or pure)—are needed. Each kind has its place; and although it is a stylistic error to mix the two indiscriminately, clearly the two have much to do with one another. However this may be, this book is a book of derivations of applied mathematics. The derivations here proceed by a purely applied approach.

Figure 1.1: Two triangles.



1.2.2 Mathematical extension

Profound results in mathematics are occasionally achieved simply by extending results already known. For example, negative integers and their properties can be discovered by counting backward—3, 2, 1, 0—then asking what follows (precedes?) 0 in the countdown and what properties this new, negative integer must have to interact smoothly with the already known positives. The astonishing Euler’s formula (§ 5.4) is discovered by a similar but more sophisticated mathematical extension.

More often, however, the results achieved by extension are unsurprising and not very interesting in themselves. Such extended results are the faithful servants of mathematical rigor. Consider for example the triangle on the left of Fig. 1.1. This triangle is evidently composed of two right triangles of areas

$$\begin{aligned} A_1 &= \frac{b_1 h}{2}, \\ A_2 &= \frac{b_2 h}{2} \end{aligned}$$

(each right triangle is exactly half a rectangle). Hence the main triangle’s area is

$$A = A_1 + A_2 = \frac{(b_1 + b_2)h}{2} = \frac{bh}{2}.$$

Very well. What about the triangle on the right? Its b_1 is not shown on the figure, and what is that $-b_2$, anyway? Answer: the triangle is composed of the *difference* of two right triangles, with b_1 the base of the larger, overall one: $b_1 = b + (-b_2)$. The b_2 is negative because the sense of the small right triangle’s area in the proof is negative: the small area is subtracted from

the large rather than added. By extension on this basis, the main triangle's area is again seen to be $A = bh/2$. The proof is exactly the same. In fact, once the central idea of adding two right triangles is grasped, the extension is really rather obvious—too obvious to be allowed to burden such a book as this.

Excepting the uncommon cases where extension reveals something interesting or new, this book generally leaves the mere extension of proofs—including the validation of edge cases and over-the-edge cases—as an exercise to the interested reader.

1.3 Complex numbers and complex variables

More than a mastery of mere logical details, it is an holistic view of the mathematics and of its use in the modeling of physical systems which is the mark of the applied mathematician. A *feel* for the math is the great thing. Formal definitions, axioms, symbolic algebras and the like, though often useful, are felt to be secondary. The book's rapidly staged development of complex numbers and complex variables is planned on this sensibility.

Sections 2.12, 3.10, 3.11, 4.3.3, 4.4, 6.2, 9.5 and 9.6.5, plus all of Chs. 5 and 8, constitute the book's principal stages of complex development. In these sections and throughout the book, the reader comes to appreciate that most mathematical properties which apply for real numbers apply equally for complex, that few properties concern real numbers alone.

Pure mathematics develops an abstract theory of the complex variable.¹ The abstract theory is quite beautiful. However, its arc takes off too late and flies too far from applications for such a book as this. Less beautiful but more practical paths to the topic exist;² this book leads the reader along one of these.

1.4 On the text

The book gives numerals in hexadecimal. It denotes variables in Greek letters as well as Roman. Readers unfamiliar with the hexadecimal notation will find a brief orientation thereto in Appendix A. Readers unfamiliar with the Greek alphabet will find it in Appendix B.

Licensed to the public under the GNU General Public Licence [14], version 2, this book meets the Debian Free Software Guidelines [9].

¹[12][39][21]

²See Ch. 8's footnote 8.

If you cite an equation, section, chapter, figure or other item from this book, it is recommended that you include in your citation the book's precise draft date as given on the title page. The reason is that equation numbers, chapter numbers and the like are numbered automatically by the \LaTeX typesetting software: such numbers can change arbitrarily from draft to draft. If an example citation helps, see [4] in the bibliography.

Chapter 2

Classical algebra and geometry

Arithmetic and the simplest elements of classical algebra and geometry, we learn as children. Few readers will want this book to begin with a treatment of $1 + 1 = 2$; or of how to solve $3x - 2 = 7$. However, there are some basic points which do seem worth touching. The book starts with these.

2.1 Basic arithmetic relationships

This section states some arithmetical rules.

2.1.1 Commutivity, associativity, distributivity, identity and inversion

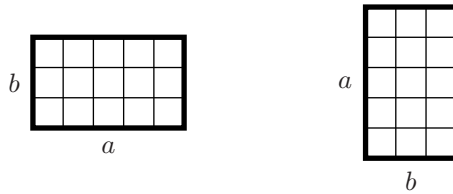
Table 2.1 lists several arithmetical rules, each of which applies not only to real numbers but equally to the complex numbers of § 2.12. Most of the rules are appreciated at once if the meaning of the symbols is understood. In the case of multiplicative commutivity, one imagines a rectangle with sides of lengths a and b , then the same rectangle turned on its side, as in Fig. 2.1: since the area of the rectangle is the same in either case, and since the area is the length times the width in either case (the area is more or less a matter of counting the little squares), evidently multiplicative commutivity holds. A similar argument validates multiplicative associativity, except that here we compute the *volume* of a three-dimensional rectangular box, which box we turn various ways.¹

¹[39, Ch. 1]

Table 2.1: Basic properties of arithmetic.

$a + b$	$=$	$b + a$	Additive commutativity
$a + (b + c)$	$=$	$(a + b) + c$	Additive associativity
$a + 0 = 0 + a$	$=$	a	Additive identity
$a + (-a)$	$=$	0	Additive inversion
ab	$=$	ba	Multiplicative commutativity
$(a)(bc)$	$=$	$(ab)(c)$	Multiplicative associativity
$(a)(1) = (1)(a)$	$=$	a	Multiplicative identity
$(a)(1/a)$	$=$	1	Multiplicative inversion
$(a)(b + c)$	$=$	$ab + ac$	Distributivity

Figure 2.1: Multiplicative commutativity.



Multiplicative inversion lacks an obvious interpretation when $a = 0$. Loosely,

$$\frac{1}{0} = \infty.$$

But since $3/0 = \infty$ also, surely either the zero or the infinity, or both, somehow differ in the latter case.

Looking ahead in the book, we note that the multiplicative properties do not always hold for more general linear transformations. For example, matrix multiplication is not commutative and vector cross-multiplication is not associative. Where associativity does not hold and parentheses do not otherwise group, right-to-left association is notationally implicit:²

$$\mathbf{A} \times \mathbf{B} \times \mathbf{C} = \mathbf{A} \times (\mathbf{B} \times \mathbf{C}).$$

The sense of it is that the thing on the left ($\mathbf{A} \times$) *operates* on the thing on the right ($\mathbf{B} \times \mathbf{C}$). (In the rare case in which the question arises, you may want to use parentheses anyway.)

2.1.2 Negative numbers

Consider that

$$\begin{aligned} (+a)(+b) &= +ab, \\ (+a)(-b) &= -ab, \\ (-a)(+b) &= -ab, \\ (-a)(-b) &= +ab. \end{aligned}$$

The first three of the four equations are unsurprising, but the last is interesting. Why would a negative count $-a$ of a negative quantity $-b$ come to

²The fine C and C++ programming languages are unfortunately stuck with the reverse order of association, along with division inharmoniously on the same level of syntactic precedence as multiplication. Standard mathematical notation is more elegant:

$$abc/uvw = \frac{(a)(bc)}{(u)(vw)}.$$

a positive product $+ab$? To see why, consider the progression

$$\begin{array}{rcl}
 & \vdots & \\
 (+3)(-b) & = & -3b, \\
 (+2)(-b) & = & -2b, \\
 (+1)(-b) & = & -1b, \\
 (0)(-b) & = & 0b, \\
 (-1)(-b) & = & +1b, \\
 (-2)(-b) & = & +2b, \\
 (-3)(-b) & = & +3b, \\
 & \vdots &
 \end{array}$$

The logic of arithmetic demands that the product of two negative numbers be positive for this reason.

2.1.3 Inequality

If³

$$a < b,$$

then necessarily

$$a + x < b + x.$$

However, the relationship between ua and ub depends on the sign of u :

$$\begin{array}{ll}
 ua < ub & \text{if } u > 0; \\
 ua > ub & \text{if } u < 0.
 \end{array}$$

Also,

$$\frac{1}{a} > \frac{1}{b}.$$

2.1.4 The change of variable

The applied mathematician very often finds it convenient *to change variables*, introducing new symbols to stand in place of old. For this we have

³Few readers attempting this book will need to be reminded that $<$ means “is less than,” that $>$ means “is greater than,” or that \leq and \geq respectively mean “is less than or equal to” and “is greater than or equal to.”

the *change of variable* or *assignment* notation⁴

$$Q \leftarrow P.$$

This means, “in place of P , put Q ”; or, “let Q now equal P .” For example, if $a^2 + b^2 = c^2$, then the change of variable $2\mu \leftarrow a$ yields the new form $(2\mu)^2 + b^2 = c^2$.

Similar to the change of variable notation is the definition notation

$$Q \equiv P.$$

This means, “let the new symbol Q represent P .”⁵

The two notations logically mean about the same thing. Subjectively, $Q \equiv P$ identifies a quantity P sufficiently interesting to be given a permanent name Q , whereas $Q \leftarrow P$ implies nothing especially interesting about P or Q ; it just introduces a (perhaps temporary) new symbol Q to ease the algebra. The concepts grow clearer as examples of the usage arise in the book.

2.2 Quadratics

Differences and sums of squares are conveniently factored as

$$\begin{aligned} a^2 - b^2 &= (a + b)(a - b), \\ a^2 + b^2 &= (a + ib)(a - ib), \\ a^2 - 2ab + b^2 &= (a - b)^2, \\ a^2 + 2ab + b^2 &= (a + b)^2 \end{aligned} \tag{2.1}$$

(where i is the *imaginary unit*, a number defined such that $i^2 = -1$, introduced in more detail in § 2.12 below). Useful as these four forms are,

⁴There appears to exist no broadly established standard mathematical notation for the change of variable, other than the $=$ equal sign, which regrettably does not fill the role well. One can indeed use the equal sign, but then what does the change of variable $k = k + 1$ mean? It looks like a claim that k and $k + 1$ are the same, which is impossible. The notation $k \leftarrow k + 1$ by contrast is unambiguous; it means to increment k by one. However, the latter notation admittedly has seen only scattered use in the literature.

The C and C++ programming languages use $==$ for equality and $=$ for assignment (change of variable), as the reader may be aware.

⁵One would never write $k \equiv k + 1$. Even $k \leftarrow k + 1$ can confuse readers inasmuch as it appears to imply two different values for the same symbol k , but the latter notation is sometimes used anyway when new symbols are unwanted or because more precise alternatives (like $k_n = k_{n-1} + 1$) seem overwrought. Still, usually it is better to introduce a new symbol, as in $j \leftarrow k + 1$.

In some books, \equiv is printed as \triangleq .

however, none of them can directly factor the more general quadratic⁶ expression

$$z^2 - 2\beta z + \gamma^2.$$

To factor this, we *complete the square*, writing

$$\begin{aligned} z^2 - 2\beta z + \gamma^2 &= z^2 - 2\beta z + \gamma^2 + (\beta^2 - \gamma^2) - (\beta^2 - \gamma^2) \\ &= z^2 - 2\beta z + \beta^2 - (\beta^2 - \gamma^2) \\ &= (z - \beta)^2 - (\beta^2 - \gamma^2). \end{aligned}$$

The expression evidently has roots⁷ where

$$(z - \beta)^2 = (\beta^2 - \gamma^2),$$

or in other words where

$$z = \beta \pm \sqrt{\beta^2 - \gamma^2}. \quad (2.2)$$

This suggests the factoring⁸

$$z^2 - 2\beta z + \gamma^2 = (z - z_1)(z - z_2), \quad (2.3)$$

where z_1 and z_2 are the two values of z given by (2.2).

It follows that the two solutions of the quadratic equation

$$z^2 = 2\beta z - \gamma^2 \quad (2.4)$$

are those given by (2.2), which is called *the quadratic formula*.⁹ (*Cubic* and *quartic formulas* also exist respectively to extract the roots of polynomials of third and fourth order, but they are much harder. See Ch. 10 and its Tables 10.1 and 10.2.)

⁶The adjective *quadratic* refers to the algebra of expressions in which no term has greater than second order. Examples of quadratic expressions include x^2 , $2x^2 - 7x + 3$ and $x^2 + 2xy + y^2$. By contrast, the expressions $x^3 - 1$ and $5x^2y$ are *cubic* not quadratic because they contain third-order terms. First-order expressions like $x + 1$ are *linear*; zeroth-order expressions like 3 are *constant*. Expressions of fourth and fifth order are *quartic* and *quintic*, respectively. (If not already clear from the context, *order* basically refers to the number of variables multiplied together in a term. The term $5x^2y = 5[x][x][y]$ is of third order, for instance.)

⁷A *root* of $f(z)$ is a value of z for which $f(z) = 0$. See § 2.11.

⁸It suggests it because the expressions on the left and right sides of (2.3) are both quadratic (the highest power is z^2) and have the same roots. Substituting into the equation the values of z_1 and z_2 and simplifying proves the suggestion correct.

⁹The form of the quadratic formula which usually appears in print is

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a},$$

2.3 Notation for series sums and products

Sums and products of series arise so frequently in mathematical work that one finds it convenient to define terse notations to express them. The summation notation

$$\sum_{k=a}^b f(k)$$

means to let k equal each of $a, a+1, a+2, \dots, b$ in turn, evaluating the function $f(k)$ at each k , then adding the several $f(k)$. For example,¹⁰

$$\sum_{k=3}^6 k^2 = 3^2 + 4^2 + 5^2 + 6^2 = 0x56.$$

The similar multiplication notation

$$\prod_{j=a}^b f(j)$$

means to *multiply* the several $f(j)$ rather than to add them. The symbols \sum and \prod come respectively from the Greek letters for S and P, and may be regarded as standing for “Sum” and “Product.” The j or k is a *dummy variable, index of summation* or *loop counter*—a variable with no independent existence, used only to facilitate the addition or multiplication of the series.¹¹ (Nothing prevents one from writing \prod_k rather than \prod_j , incidentally. For a dummy variable, one can use any letter one likes. However, the general habit of writing \sum_k and \prod_j proves convenient at least in § 4.5.2 and Ch. 8, so we start now.)

which solves the quadratic $ax^2 + bx + c = 0$. However, this writer finds the form (2.2) easier to remember. For example, by (2.2) in light of (2.4), the quadratic

$$z^2 = 3z - 2$$

has the solutions

$$z = \frac{3}{2} \pm \sqrt{\left(\frac{3}{2}\right)^2 - 2} = 1 \text{ or } 2.$$

¹⁰The hexadecimal numeral 0x56 represents the same number the decimal numeral 86 represents. The book’s preface explains why the book represents such numbers in hexadecimal. Appendix A tells how to read the numerals.

¹¹Section 7.3 speaks further of the dummy variable.

The product shorthand

$$\begin{aligned} n! &\equiv \prod_{j=1}^n j, \\ n!/m! &\equiv \prod_{j=m+1}^n j, \end{aligned}$$

is very frequently used. The notation $n!$ is pronounced “ n factorial.” Regarding the notation $n!/m!$, this can of course be regarded correctly as $n!$ divided by $m!$, but it usually proves more amenable to regard the notation as a single unit.¹²

Because multiplication in its more general sense as linear transformation (§ 11.1.1) is not always commutative, we specify that

$$\prod_{j=a}^b f(j) = [f(b)][f(b-1)][f(b-2)] \cdots [f(a+2)][f(a+1)][f(a)]$$

rather than the reverse order of multiplication.¹³ Multiplication proceeds from right to left. In the event that the reverse order of multiplication is needed, we shall use the notation

$$\prod_{j=a}^b f(j) = [f(a)][f(a+1)][f(a+2)] \cdots [f(b-2)][f(b-1)][f(b)].$$

Note that for the sake of definitional consistency,

$$\begin{aligned} \sum_{k=N+1}^N f(k) &= 0 + \sum_{k=N+1}^N f(k) = 0, \\ \prod_{j=N+1}^N f(j) &= (1) \prod_{j=N+1}^N f(j) = 1. \end{aligned}$$

This means among other things that

$$0! = 1. \tag{2.5}$$

¹²One reason among others for this is that factorials rapidly multiply to extremely large sizes, overflowing computer registers during numerical computation. If you can avoid unnecessary multiplication by regarding $n!/m!$ as a single unit, this is a win.

¹³The extant mathematical literature lacks an established standard on the order of multiplication implied by the “ \prod ” symbol, but this is the order we shall use in this book.

On first encounter, such \sum and \prod notation seems a bit overwrought. Admittedly it is easier for the beginner to read “ $f(1) + f(2) + \cdots + f(N)$ ” than “ $\sum_{k=1}^N f(k)$.” However, experience shows the latter notation to be extremely useful in expressing more sophisticated mathematical ideas. We shall use such notation extensively in this book.

2.4 The arithmetic series

A simple yet useful application of the series sum of § 2.3 is the *arithmetic series*

$$\sum_{k=a}^b k = a + (a+1) + (a+2) + \cdots + b.$$

Pairing a with b , then $a+1$ with $b-1$, then $a+2$ with $b-2$, etc., the average of each pair is $[a+b]/2$; thus the average of the entire series is $[a+b]/2$. (The pairing may or may not leave an unpaired element at the series midpoint $k = [a+b]/2$, but this changes nothing.) The series has $b - a + 1$ terms. Hence,

$$\sum_{k=a}^b k = (b - a + 1) \frac{a + b}{2}. \quad (2.6)$$

Success with this arithmetic series leads one to wonder about the *geometric series* $\sum_{k=0}^{\infty} z^k$. Section 2.6.4 addresses that point.

2.5 Powers and roots

This necessarily tedious section discusses powers and roots. It offers no surprises. Table 2.2 summarizes its definitions and results. Readers seeking more rewarding reading may prefer just to glance at the table then to skip directly to the start of the next section.

In this section, the exponents k, m, n, p, q, r and s are integers,¹⁴ but the exponents a and b are arbitrary real numbers.

¹⁴In case the word is unfamiliar to the reader who has learned arithmetic in another language than English, the *integers* are the negative, zero and positive counting numbers $\dots, -3, -2, -1, 0, 1, 2, 3, \dots$. The corresponding adjective is *integral* (although the word “integral” is also used as a noun and adjective indicating an infinite sum of infinitesimals; see Ch. 7). Traditionally, the letters i, j, k, m, n, M and N are used to represent integers (i is sometimes avoided because the same letter represents the imaginary unit), but this section needs more integer symbols so it uses p, q, r and s , too.

Table 2.2: Power properties and definitions.

$$\begin{aligned}
z^n &\equiv \prod_{j=1}^n z, \quad n \geq 0 \\
z &= (z^{1/n})^n = (z^n)^{1/n} \\
\sqrt{z} &\equiv z^{1/2} \\
(uv)^a &= u^a v^a \\
z^{p/q} &= (z^{1/q})^p = (z^p)^{1/q} \\
z^{ab} &= (z^a)^b = (z^b)^a \\
z^{a+b} &= z^a z^b \\
z^{a-b} &= \frac{z^a}{z^b} \\
z^{-b} &= \frac{1}{z^b}
\end{aligned}$$

2.5.1 Notation and integral powers

The power notation

$$z^n$$

indicates the number z , multiplied by itself n times. More formally, when the *exponent* n is a nonnegative integer,¹⁵

$$z^n \equiv \prod_{j=1}^n z. \tag{2.7}$$

¹⁵The symbol “ \equiv ” means “ $=$ ”, but it further usually indicates that the expression on its right serves to define the expression on its left. Refer to § 2.1.4.

For example,¹⁶

$$\begin{aligned} z^3 &= (z)(z)(z), \\ z^2 &= (z)(z), \\ z^1 &= z, \\ z^0 &= 1. \end{aligned}$$

Notice that in general,

$$z^{n-1} = \frac{z^n}{z}.$$

This leads us to extend the definition to negative integral powers with

$$z^{-n} = \frac{1}{z^n}. \quad (2.8)$$

From the foregoing it is plain that

$$\begin{aligned} z^{m+n} &= z^m z^n, \\ z^{m-n} &= \frac{z^m}{z^n}, \end{aligned} \quad (2.9)$$

for any integral m and n . For similar reasons,

$$z^{mn} = (z^m)^n = (z^n)^m. \quad (2.10)$$

On the other hand from multiplicative associativity and commutivity,

$$(uv)^n = u^n v^n. \quad (2.11)$$

2.5.2 Roots

Fractional powers are not something we have defined yet, so for consistency with (2.10) we let

$$(u^{1/n})^n = u.$$

This has $u^{1/n}$ as the number which, raised to the n th power, yields u . Setting

$$v = u^{1/n},$$

¹⁶The case 0^0 is interesting because it lacks an obvious interpretation. The specific interpretation depends on the nature and meaning of the two zeros. For interest, if $E \equiv 1/\epsilon$, then

$$\lim_{\epsilon \rightarrow 0^+} \epsilon^\epsilon = \lim_{E \rightarrow \infty} \left(\frac{1}{E} \right)^{1/E} = \lim_{E \rightarrow \infty} E^{-1/E} = \lim_{E \rightarrow \infty} e^{-(\ln E)/E} = e^0 = 1.$$

it follows by successive steps that

$$\begin{aligned} v^n &= u, \\ (v^n)^{1/n} &= u^{1/n}, \\ (v^n)^{1/n} &= v. \end{aligned}$$

Taking the u and v formulas together, then,

$$(z^{1/n})^n = z = (z^n)^{1/n} \quad (2.12)$$

for any z and integral n .

The number $z^{1/n}$ is called the n th root of z —or in the very common case $n = 2$, the *square root* of z , often written

$$\sqrt{z}.$$

When z is real and nonnegative, the last notation is usually implicitly taken to mean the real, nonnegative square root. In any case, the power and root operations mutually invert one another.

What about powers expressible neither as n nor as $1/n$, such as the $3/2$ power? If z and w are numbers related by

$$w^q = z,$$

then

$$w^{pq} = z^p.$$

Taking the q th root,

$$w^p = (z^p)^{1/q}.$$

But $w = z^{1/q}$, so this is

$$(z^{1/q})^p = (z^p)^{1/q},$$

which says that it does not matter whether one applies the power or the root first; the result is the same. Extending (2.10) therefore, we define $z^{p/q}$ such that

$$(z^{1/q})^p = z^{p/q} = (z^p)^{1/q}. \quad (2.13)$$

Since any real number can be approximated arbitrarily closely by a ratio of integers, (2.13) implies a power definition for all real exponents.

Equation (2.13) is this subsection's main result. However, § 2.5.3 will find it useful if we can also show here that

$$(z^{1/q})^{1/s} = z^{1/qs} = (z^{1/s})^{1/q}. \quad (2.14)$$

The proof is straightforward. If

$$w \equiv z^{1/qs},$$

then raising to the qs power yields

$$(w^s)^q = z.$$

Successively taking the q th and s th roots gives

$$w = (z^{1/q})^{1/s}.$$

By identical reasoning,

$$w = (z^{1/s})^{1/q}.$$

But since $w \equiv z^{1/qs}$, the last two equations imply (2.14), as we have sought.

2.5.3 Powers of products and powers of powers

Per (2.11),

$$(uv)^p = u^p v^p.$$

Raising this equation to the $1/q$ power, we have that

$$\begin{aligned} (uv)^{p/q} &= [u^p v^p]^{1/q} \\ &= \left[(u^p)^{q/q} (v^p)^{q/q} \right]^{1/q} \\ &= \left[(u^{p/q})^q (v^{p/q})^q \right]^{1/q} \\ &= \left[(u^{p/q})(v^{p/q}) \right]^{q/q} \\ &= u^{p/q} v^{p/q}. \end{aligned}$$

In other words

$$(uv)^a = u^a v^a \tag{2.15}$$

for any real a .

On the other hand, per (2.10),

$$z^{pr} = (z^p)^r.$$

Raising this equation to the $1/qs$ power and applying (2.10), (2.13) and (2.14) to reorder the powers, we have that

$$z^{(p/q)(r/s)} = (z^{p/q})^{r/s}.$$

By identical reasoning,

$$z^{(p/q)(r/s)} = (z^{r/s})^{p/q}.$$

Since p/q and r/s can approximate any real numbers with arbitrary precision, this implies that

$$(z^a)^b = z^{ab} = (z^b)^a \quad (2.16)$$

for any real a and b .

2.5.4 Sums of powers

With (2.9), (2.15) and (2.16), one can reason that

$$z^{(p/q)+(r/s)} = (z^{ps+rq})^{1/qs} = (z^{ps}z^{rq})^{1/qs} = z^{p/q}z^{r/s},$$

or in other words that

$$z^{a+b} = z^a z^b. \quad (2.17)$$

In the case that $a = -b$,

$$1 = z^{-b+b} = z^{-b} z^b,$$

which implies that

$$z^{-b} = \frac{1}{z^b}. \quad (2.18)$$

But then replacing $-b \leftarrow b$ in (2.17) leads to

$$z^{a-b} = z^a z^{-b},$$

which according to (2.18) is

$$z^{a-b} = \frac{z^a}{z^b}. \quad (2.19)$$

2.5.5 Summary and remarks

Table 2.2 on page 16 summarizes the section's definitions and results.

Looking ahead to § 2.12, § 3.11 and Ch. 5, we observe that nothing in the foregoing analysis requires the base variables z , w , u and v to be real numbers; if complex (§ 2.12), the formulas remain valid. Still, the analysis does imply that the various exponents m , n , p/q , a , b and so on are real numbers. This restriction, we shall remove later, purposely defining the action of a complex exponent to comport with the results found here. With such a definition the results apply not only for all bases but also for all exponents, real or complex.

2.6 Multiplying and dividing power series

A *power series*¹⁷ is a weighted sum of integral powers:

$$A(z) = \sum_{k=-\infty}^{\infty} a_k z^k, \quad (2.20)$$

where the several a_k are arbitrary constants. This section discusses the multiplication and division of power series.

¹⁷Another name for the *power series* is *polynomial*. The word “polynomial” usually connotes a power series with a finite number of terms, but the two names in fact refer to essentially the same thing.

Professional mathematicians use the terms more precisely. Equation (2.20), they call a “power series” only if $a_k = 0$ for all $k < 0$ —in other words, technically, not if it includes negative powers of z . They call it a “polynomial” only if it is a “power series” with a finite number of terms. They call (2.20) in general a *Laurent series*.

The name “Laurent series” is a name we shall meet again in § 8.14. In the meantime however we admit that the professionals have vaguely daunted us by adding to the name some pretty sophisticated connotations, to the point that we applied mathematicians (at least in the author’s country) seem to feel somehow unlicensed actually to use the name. We tend to call (2.20) a “power series with negative powers,” or just “a power series.”

This book follows the last usage. You however can call (2.20) a *Laurent series* if you prefer (and if you pronounce it right: “lor-ON”). That is after all exactly what it is. Nevertheless if you do use the name “Laurent series,” be prepared for people subjectively—for no particular reason—to expect you to establish complex radii of convergence, to sketch some annulus in the Argand plane, and/or to engage in other maybe unnecessary formalities. If that is not what you seek, then you may find it better just to call the thing by the less lofty name of “power series”—or better, if it has a finite number of terms, by the even humbler name of “polynomial.”

Semantics. All these names mean about the same thing, but one is expected most carefully always to give the right name in the right place. What a bother! (Someone once told the writer that the Japanese language can give different names to the same object, depending on whether the *speaker* is male or female. The power-series terminology seems to share a spirit of that kin.) If you seek just one word for the thing, the writer recommends that you call it a “power series” and then not worry too much about it until someone objects. When someone does object, you can snow him with the big word “Laurent series,” instead.

The experienced scientist or engineer may notice that the above vocabulary omits the name “Taylor series.” The vocabulary omits the name because that name fortunately remains unconfused in usage—it means quite specifically a power series without negative powers and tends to connote a representation of some particular function of interest—as we shall see in Ch. 8.

2.6.1 Multiplying power series

Given two power series

$$\begin{aligned} A(z) &= \sum_{k=-\infty}^{\infty} a_k z^k, \\ B(z) &= \sum_{k=-\infty}^{\infty} b_k z^k, \end{aligned} \tag{2.21}$$

the product of the two series is evidently

$$P(z) \equiv A(z)B(z) = \sum_{k=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} a_j b_{k-j} z^k. \tag{2.22}$$

2.6.2 Dividing power series

The quotient $Q(z) = B(z)/A(z)$ of two power series is a little harder to calculate, and there are at least two ways to do it. Section 2.6.3 below will do it by matching coefficients, but this subsection does it by long division. For example,

$$\begin{aligned} \frac{2z^2 - 3z + 3}{z - 2} &= \frac{2z^2 - 4z}{z - 2} + \frac{z + 3}{z - 2} = 2z + \frac{z + 3}{z - 2} \\ &= 2z + \frac{z - 2}{z - 2} + \frac{5}{z - 2} = 2z + 1 + \frac{5}{z - 2}. \end{aligned}$$

The strategy is to take the dividend¹⁸ $B(z)$ piece by piece, purposely choosing pieces easily divided by $A(z)$.

If you feel that you understand the example, then that is really all there is to it, and you can skip the rest of the subsection if you like. One sometimes wants to express the long division of power series more formally, however. That is what the rest of the subsection is about. (Be advised however that the cleverer technique of § 2.6.3, though less direct, is often easier and faster.)

Formally, we prepare the long division $B(z)/A(z)$ by writing

$$B(z) = A(z)Q_n(z) + R_n(z), \tag{2.23}$$

¹⁸If $Q(z)$ is a *quotient* and $R(z)$ a *remainder*, then $B(z)$ is a *dividend* (or *numerator*) and $A(z)$ a *divisor* (or *denominator*). Such are the Latin-derived names of the parts of a long division.

where $R_n(z)$ is a *remainder* (being the part of $B[z]$ *remaining* to be divided); and

$$\begin{aligned}
A(z) &= \sum_{k=-\infty}^K a_k z^k, \quad a_K \neq 0, \\
B(z) &= \sum_{k=-\infty}^N b_k z^k, \\
R_N(z) &= B(z), \\
Q_N(z) &= 0, \\
R_n(z) &= \sum_{k=-\infty}^n r_{nk} z^k, \\
Q_n(z) &= \sum_{k=n-K+1}^{N-K} q_k z^k,
\end{aligned} \tag{2.24}$$

where K and N identify the greatest orders k of z^k present in $A(z)$ and $B(z)$, respectively.

Well, that is a lot of symbology. What does it mean? The key to understanding it lies in understanding (2.23), which is not one but several equations—one equation for each value of n , where $n = N, N-1, N-2, \dots$. The dividend $B(z)$ and the divisor $A(z)$ stay the same from one n to the next, but the quotient $Q_n(z)$ and the remainder $R_n(z)$ change. At start, $Q_N(z) = 0$ while $R_N(z) = B(z)$, but the thrust of the long division process is to build $Q_n(z)$ up by wearing $R_n(z)$ down. The goal is to grind $R_n(z)$ away to nothing, to make it disappear as $n \rightarrow -\infty$.

As in the example, we pursue the goal by choosing from $R_n(z)$ an easily divisible piece containing the whole high-order term of $R_n(z)$. The piece we choose is $(r_{nn}/a_K)z^{n-K}A(z)$, which we add and subtract from (2.23) to obtain

$$B(z) = A(z) \left[Q_n(z) + \frac{r_{nn}}{a_K} z^{n-K} \right] + \left[R_n(z) - \frac{r_{nn}}{a_K} z^{n-K} A(z) \right].$$

Matching this equation against the desired iterate

$$B(z) = A(z)Q_{n-1}(z) + R_{n-1}(z)$$

and observing from the definition of $Q_n(z)$ that $Q_{n-1}(z) = Q_n(z) +$

$q_{n-K}z^{n-K}$, we find that

$$\begin{aligned} q_{n-K} &= \frac{r_{nn}}{a_K}, \\ R_{n-1}(z) &= R_n(z) - q_{n-K}z^{n-K}A(z), \end{aligned} \tag{2.25}$$

where no term remains in $R_{n-1}(z)$ higher than a z^{n-1} term.

To begin the actual long division, we initialize

$$R_N(z) = B(z),$$

for which (2.23) is trivially true. Then we iterate per (2.25) as many times as desired. If an infinite number of times, then so long as $R_n(z)$ tends to vanish as $n \rightarrow -\infty$, it follows from (2.23) that

$$\frac{B(z)}{A(z)} = Q_{-\infty}(z). \tag{2.26}$$

Iterating only a finite number of times leaves a remainder,

$$\frac{B(z)}{A(z)} = Q_n(z) + \frac{R_n(z)}{A(z)}, \tag{2.27}$$

except that it may happen that $R_n(z) = 0$ for sufficiently small n .

Table 2.3 summarizes the long-division procedure.¹⁹ In its q_{n-K} equation, the table includes also the result of § 2.6.3 below.

It should be observed in light of Table 2.3 that if²⁰

$$\begin{aligned} A(z) &= \sum_{k=K_o}^K a_k z^k, \\ B(z) &= \sum_{k=N_o}^N b_k z^k, \end{aligned}$$

then

$$R_n(z) = \sum_{k=n-(K-K_o)+1}^n r_{nk} z^k \quad \text{for all } n < N_o + (K - K_o). \tag{2.28}$$

¹⁹[41, § 3.2]

²⁰The notations K_o , a_k and z^k are usually pronounced, respectively, as “ K naught,” “ a sub k ” and “ z to the k ” (or, more fully, “ z to the k th power”)—at least in the author’s country.

Table 2.3: Dividing power series through successively smaller powers.

$$\begin{aligned}
B(z) &= A(z)Q_n(z) + R_n(z) \\
A(z) &= \sum_{k=-\infty}^K a_k z^k, \quad a_K \neq 0 \\
B(z) &= \sum_{k=-\infty}^N b_k z^k \\
R_N(z) &= B(z) \\
Q_N(z) &= 0 \\
R_n(z) &= \sum_{k=-\infty}^n r_{nk} z^k \\
Q_n(z) &= \sum_{k=n-K+1}^{N-K} q_k z^k \\
q_{n-K} &= \frac{r_{nn}}{a_K} = \frac{1}{a_K} \left(b_n - \sum_{k=n-K+1}^{N-K} a_{n-k} q_k \right) \\
R_{n-1}(z) &= R_n(z) - q_{n-K} z^{n-K} A(z) \\
\frac{B(z)}{A(z)} &= Q_{-\infty}(z)
\end{aligned}$$

That is, the remainder has order one less than the divisor has. The reason for this, of course, is that we have strategically planned the long-division iteration precisely to cause the leading term of the divisor to cancel the leading term of the remainder at each step.²¹

The long-division procedure of Table 2.3 extends the quotient $Q_n(z)$ through successively smaller powers of z . Often, however, one prefers to extend the quotient through successively *larger* powers of z , where a z^K term is $A(z)$'s term of *least* order. In this case, the long division goes by the complementary rules of Table 2.4.

2.6.3 Dividing power series by matching coefficients

There is another, sometimes quicker way to divide power series than by the long division of § 2.6.2. One can divide them by matching coefficients.²² If

$$Q_\infty(z) = \frac{B(z)}{A(z)}, \quad (2.29)$$

where

$$\begin{aligned} A(z) &= \sum_{k=K}^{\infty} a_k z^k, \quad a_K \neq 0, \\ B(z) &= \sum_{k=N}^{\infty} b_k z^k \end{aligned}$$

are known and

$$Q_\infty(z) = \sum_{k=N-K}^{\infty} q_k z^k$$

²¹If a more formal demonstration of (2.28) is wanted, then consider per (2.25) that

$$R_{m-1}(z) = R_m(z) - \frac{r_{mm}}{a_K} z^{m-K} A(z).$$

If the least-order term of $R_m(z)$ is a z^{N_o} term (as clearly is the case at least for the initial remainder $R_N[z] = B[z]$), then according to the equation so also must the least-order term of $R_{m-1}(z)$ be a z^{N_o} term, unless an even lower-order term be contributed by the product $z^{m-K} A(z)$. But that very product's term of least order is a $z^{m-(K-K_o)}$ term. Under these conditions, evidently the least-order term of $R_{m-1}(z)$ is a $z^{m-(K-K_o)}$ term when $m - (K - K_o) \leq N_o$; otherwise a z^{N_o} term. This is better stated after the change of variable $n+1 \leftarrow m$: the least-order term of $R_n(z)$ is a $z^{n-(K-K_o)+1}$ term when $n < N_o + (K - K_o)$; otherwise a z^{N_o} term.

The greatest-order term of $R_n(z)$ is by definition a z^n term. So, in summary, when $n < N_o + (K - K_o)$, the terms of $R_n(z)$ run from $z^{n-(K-K_o)+1}$ through z^n , which is exactly the claim (2.28) makes.

²²[25][12, § 2.5]

Table 2.4: Dividing power series through successively larger powers.

$$\begin{aligned}
B(z) &= A(z)Q_n(z) + R_n(z) \\
A(z) &= \sum_{k=K}^{\infty} a_k z^k, \quad a_K \neq 0 \\
B(z) &= \sum_{k=N}^{\infty} b_k z^k \\
R_N(z) &= B(z) \\
Q_N(z) &= 0 \\
R_n(z) &= \sum_{k=n}^{\infty} r_{nk} z^k \\
Q_n(z) &= \sum_{k=N-K}^{n-K-1} q_k z^k \\
q_{n-K} &= \frac{r_{nn}}{a_K} = \frac{1}{a_K} \left(b_n - \sum_{k=N-K}^{n-K-1} a_{n-k} q_k \right) \\
R_{n+1}(z) &= R_n(z) - q_{n-K} z^{n-K} A(z) \\
\frac{B(z)}{A(z)} &= Q_{\infty}(z)
\end{aligned}$$

is to be calculated, then one can multiply (2.29) through by $A(z)$ to obtain the form

$$A(z)Q_\infty(z) = B(z).$$

Expanding the left side according to (2.22) and changing the index $n \leftarrow k$ on the right side,

$$\sum_{n=N}^{\infty} \sum_{k=N-K}^{n-K} a_{n-k} q_k z^n = \sum_{n=N}^{\infty} b_n z^n.$$

But for this to hold for all z , the coefficients must match for each n :

$$\sum_{k=N-K}^{n-K} a_{n-k} q_k = b_n, \quad n \geq N.$$

Transferring all terms but $a_K q_{n-K}$ to the equation's right side and dividing by a_K , we have that

$$q_{n-K} = \frac{1}{a_K} \left(b_n - \sum_{k=N-K}^{n-K-1} a_{n-k} q_k \right), \quad n \geq N. \quad (2.30)$$

Equation (2.30) computes the coefficients of $Q(z)$, each coefficient depending on the coefficients earlier computed.

The coefficient-matching technique of this subsection is easily adapted to the division of series in decreasing, rather than increasing, powers of z if needed or desired. The adaptation is left as an exercise to the interested reader, but Tables 2.3 and 2.4 incorporate the technique both ways.

Admittedly, the fact that (2.30) yields a sequence of coefficients does not necessarily mean that the resulting power series $Q_\infty(z)$ converges to some definite value over a given domain. Consider for instance (2.34), which diverges when $|z| > 1$, even though all its coefficients are known. At least (2.30) is correct when $Q_\infty(z)$ does converge. Even when $Q_\infty(z)$ as such does not converge, however, often what interest us are only the series' first several terms

$$Q_n(z) = \sum_{k=N-K}^{n-K-1} q_k z^k.$$

In this case,

$$Q_\infty(z) = \frac{B(z)}{A(z)} = Q_n(z) + \frac{R_n(z)}{A(z)} \quad (2.31)$$

and convergence is not an issue. Solving (2.31) for $R_n(z)$,

$$R_n(z) = B(z) - A(z)Q_n(z). \quad (2.32)$$

2.6.4 Common power-series quotients and the geometric series

Frequently encountered power-series quotients, calculated by the long division of § 2.6.2, computed by the coefficient matching of § 2.6.3, and/or verified by multiplying, include²³

$$\frac{1}{1 \pm z} = \begin{cases} \sum_{k=0}^{\infty} (\mp)^k z^k, & |z| < 1; \\ -\sum_{k=-\infty}^{-1} (\mp)^k z^k, & |z| > 1. \end{cases} \quad (2.33)$$

Equation (2.33) almost incidentally answers a question which has arisen in § 2.4 and which often arises in practice: to what total does the infinite *geometric series* $\sum_{k=0}^{\infty} z^k$, $|z| < 1$, sum? Answer: it sums exactly to $1/(1 - z)$. However, there is a simpler, more aesthetic way to demonstrate the same thing, as follows. Let

$$S \equiv \sum_{k=0}^{\infty} z^k, \quad |z| < 1.$$

Multiplying by z yields

$$zS \equiv \sum_{k=1}^{\infty} z^k.$$

Subtracting the latter equation from the former leaves

$$(1 - z)S = 1,$$

which, after dividing by $1 - z$, implies that

$$S \equiv \sum_{k=0}^{\infty} z^k = \frac{1}{1 - z}, \quad |z| < 1, \quad (2.34)$$

as was to be demonstrated.

2.6.5 Variations on the geometric series

Besides being more aesthetic than the long division of § 2.6.2, the difference technique of § 2.6.4 permits one to extend the basic geometric series in

²³The notation $|z|$ represents the magnitude of z . For example, $|5| = 5$, but also $|-5| = 5$.

several ways. For instance, the sum

$$S_1 \equiv \sum_{k=0}^{\infty} k z^k, \quad |z| < 1$$

(which arises in, among others, Planck's quantum blackbody radiation calculation²⁴), we can compute as follows. We multiply the unknown S_1 by z , producing

$$z S_1 = \sum_{k=0}^{\infty} k z^{k+1} = \sum_{k=1}^{\infty} (k-1) z^k.$$

We then subtract $z S_1$ from S_1 , leaving

$$(1-z) S_1 = \sum_{k=0}^{\infty} k z^k - \sum_{k=1}^{\infty} (k-1) z^k = \sum_{k=1}^{\infty} z^k = z \sum_{k=0}^{\infty} z^k = \frac{z}{1-z},$$

where we have used (2.34) to collapse the last sum. Dividing by $1-z$, we arrive at

$$S_1 \equiv \sum_{k=0}^{\infty} k z^k = \frac{z}{(1-z)^2}, \quad |z| < 1, \quad (2.35)$$

which was to be found.

Further series of the kind, such as $\sum_k k^2 z^k$, $\sum_k (k+1)(k) z^k$, $\sum_k k^3 z^k$, etc., can be calculated in like manner as the need for them arises.

2.7 Indeterminate constants, independent variables and dependent variables

Mathematical models use *indeterminate constants*, *independent variables* and *dependent variables*. The three are best illustrated by example. Consider the time t a sound needs to travel from its source to a distant listener:

$$t = \frac{\Delta r}{v_{\text{sound}}},$$

where Δr is the distance from source to listener and v_{sound} is the speed of sound. Here, v_{sound} is an indeterminate constant (given particular atmospheric conditions, it doesn't vary), Δr is an independent variable, and t is a dependent variable. The model gives t as a function of Δr ; so, if you tell the model how far the listener sits from the sound source, the model

²⁴[29]

returns the time needed for the sound to propagate from one to the other. Note that the abstract validity of the model does not necessarily depend on whether we actually know the right figure for v_{sound} (if I tell you that sound goes at 500 m/s, but later you find out that the real figure is 331 m/s, it probably doesn't ruin the theoretical part of your analysis; you just have to recalculate numerically). Knowing the figure is not the point. The point is that conceptually there preëxists some right figure for the indeterminate constant; that sound goes at some constant speed—whatever it is—and that we can calculate the delay in terms of this.

Although there exists a definite philosophical distinction between the three kinds of quantity, nevertheless it cannot be denied that which particular quantity is an indeterminate constant, an independent variable or a dependent variable often depends upon one's immediate point of view. The same model in the example would remain valid if atmospheric conditions were changing (v_{sound} would then be an independent variable) or if the model were used in designing a musical concert hall²⁵ to suffer a maximum acceptable sound time lag from the stage to the hall's back row (t would then be an independent variable; Δr , dependent). Occasionally we go so far as deliberately to change our point of view in mid-analysis, now regarding as an independent variable what a moment ago we had regarded as an indeterminate constant, for instance (a typical case of this arises in the solution of differential equations by the method of unknown coefficients, § 9.4). Such a shift of viewpoint is fine, so long as we remember that there is a difference between the three kinds of quantity and we keep track of which quantity is which kind to us at the moment.

The main reason it matters which symbol represents which of the three kinds of quantity is that in calculus, one analyzes how change in independent variables affects dependent variables as indeterminate constants remain

²⁵Math books are funny about examples like this. Such examples remind one of the kind of calculation one encounters in a childhood arithmetic textbook, as of the quantity of air contained in an astronaut's round helmet. One could calculate the quantity of water in a kitchen mixing bowl just as well, but astronauts' helmets are so much more interesting than bowls, you see.

The chance that the typical reader will ever specify the dimensions of a real musical concert hall is of course vanishingly small. However, it is the idea of the example that matters here, because the chance that the typical reader will ever specify *something* technical is quite large. Although sophisticated models with many factors and terms do indeed play a major role in engineering, the great majority of practical engineering calculations—for quick, day-to-day decisions where small sums of money and negligible risk to life are at stake—are done with models hardly more sophisticated than the one shown here. So, maybe the concert-hall example is not so unreasonable, after all.

fixed.

(Section 2.3 has introduced the dummy variable, which the present section's threefold taxonomy seems to exclude. However, in fact, most dummy variables are just independent variables—a few are dependent variables—whose scope is restricted to a particular expression. Such a dummy variable does not seem very “independent,” of course; but its dependence is on the operator controlling the expression, not on some other variable within the expression. Within the expression, the dummy variable fills the role of an independent variable; without, it fills no role because logically it does not exist there. Refer to §§ 2.3 and 7.3.)

2.8 Exponentials and logarithms

In § 2.5 we have considered the power operation z^a , where (in § 2.7's language) the independent variable z is the base and the indeterminate constant a is the exponent. There is another way to view the power operation, however. One can view it as the *exponential* operation

$$a^z,$$

where the variable z is the exponent and the constant a is the base.

2.8.1 The logarithm

The exponential operation follows the same laws the power operation follows, but because the variable of interest is now the exponent rather than the base, the inverse operation is not the root but rather the *logarithm*:

$$\log_a(a^z) = z. \quad (2.36)$$

The logarithm $\log_a w$ answers the question, “What power must I raise a to, to get w ?”

Raising a to the power of the last equation, we have that

$$a^{\log_a(a^z)} = a^z.$$

With the change of variable $w \leftarrow a^z$, this is

$$a^{\log_a w} = w. \quad (2.37)$$

Hence, the exponential and logarithmic operations mutually invert one another.

2.8.2 Properties of the logarithm

The basic properties of the logarithm include that

$$\log_a uv = \log_a u + \log_a v, \quad (2.38)$$

$$\log_a \frac{u}{v} = \log_a u - \log_a v, \quad (2.39)$$

$$\log_a (w^z) = z \log_a w, \quad (2.40)$$

$$w^z = a^{z \log_a w}, \quad (2.41)$$

$$\log_b w = \frac{\log_a w}{\log_a b}. \quad (2.42)$$

Of these, (2.38) follows from the steps

$$\begin{aligned} (uv) &= (u)(v), \\ (a^{\log_a uv}) &= (a^{\log_a u})(a^{\log_a v}), \\ a^{\log_a uv} &= a^{\log_a u + \log_a v}, \end{aligned}$$

and (2.39) follows by similar reasoning. Equations (2.40) and (2.41) follow from the steps

$$\begin{aligned} w^z &= (w^z) = (w)^z, \\ w^z &= a^{\log_a (w^z)} = (a^{\log_a w})^z, \\ w^z &= a^{\log_a (w^z)} = a^{z \log_a w}. \end{aligned}$$

Equation (2.42) follows from the steps

$$\begin{aligned} w &= b^{\log_b w}, \\ \log_a w &= \log_a (b^{\log_b w}), \\ \log_a w &= \log_b w \log_a b. \end{aligned}$$

Among other purposes, (2.38) through (2.42) serve respectively to transform products to sums, quotients to differences, powers to products, exponentials to differently based exponentials, and logarithms to differently based logarithms. Table 2.5 repeats the equations along with (2.36) and (2.37) (which also emerge as restricted forms of eqns. 2.40 and 2.41), thus summarizing the general properties of the logarithm.

2.9 Triangles and other polygons: simple facts

This section gives simple facts about triangles and other polygons.

Table 2.5: General properties of the logarithm.

$$\begin{aligned}
\log_a uv &= \log_a u + \log_a v \\
\log_a \frac{u}{v} &= \log_a u - \log_a v \\
\log_a(w^z) &= z \log_a w \\
w^z &= a^{z \log_a w} \\
\log_b w &= \frac{\log_a w}{\log_a b} \\
\log_a(a^z) &= z \\
w &= a^{\log_a w}
\end{aligned}$$

2.9.1 Triangle area

The area of a right triangle²⁶ is half the area of the corresponding rectangle. This is seen by splitting a rectangle down its diagonal into a pair of right triangles of equal size. The fact that *any* triangle's area is half its base length times its height is seen by dropping a perpendicular from one point of the triangle to the opposite side (see Fig. 1.1 on page 4), dividing the triangle into two right triangles, for each of which the fact is true. In algebraic symbols,

$$A = \frac{bh}{2}, \quad (2.43)$$

where A stands for area, b for base length, and h for perpendicular height.

2.9.2 The triangle inequalities

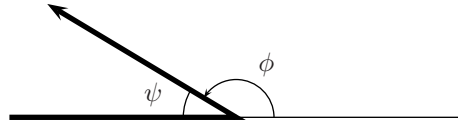
Any two sides of a triangle together are longer than the third alone, which itself is longer than the difference between the two. In symbols,

$$|a - b| < c < a + b, \quad (2.44)$$

where a , b and c are the lengths of a triangle's three sides. These are the *triangle inequalities*. The truth of the sum inequality $c < a + b$, is seen by sketching some triangle on a sheet of paper and asking: if c is the direct route between two points and $a + b$ is an indirect route, then how can $a + b$ not be longer? Of course the sum inequality is equally good on any of the

²⁶A *right triangle* is a triangle, one of whose three angles is perfectly square.

Figure 2.2: The sum of a triangle's inner angles: turning at the corner.



triangle's three sides, so one can write $a < c + b$ and $b < c + a$ just as well as $c < a + b$. Rearranging the a and b inequalities, we have that $a - b < c$ and $b - a < c$, which together say that $|a - b| < c$. The last is the difference inequality, completing (2.44)'s proof.

2.9.3 The sum of interior angles

A triangle's three interior angles²⁷ sum to $2\pi/2$. One way to see the truth of this fact is to imagine a small car rolling along one of the triangle's sides. Reaching the corner, the car turns to travel along the next side, and so on round all three corners to complete a circuit, returning to the start. Since the car again faces the original direction, we reason that it has turned a total of 2π : a full revolution. But the angle ϕ the car turns at a corner and the triangle's inner angle ψ there together form the straight angle $2\pi/2$ (the sharper the inner angle, the more the car turns: see Fig. 2.2). In mathematical notation,

$$\begin{aligned}\phi_1 + \phi_2 + \phi_3 &= 2\pi, \\ \phi_k + \psi_k &= \frac{2\pi}{2}, \quad k = 1, 2, 3,\end{aligned}$$

where ψ_k and ϕ_k are respectively the triangle's inner angles and the angles through which the car turns. Solving the latter equations for ϕ_k and

²⁷Many or most readers will already know the notation 2π and its meaning as the angle of full revolution. For those who do not, the notation is introduced more properly in §§ 3.1, 3.6 and 8.11 below. Briefly, however, the symbol 2π represents a complete turn, a full circle, a spin to face the same direction as before. Hence $2\pi/4$ represents a square turn or right angle.

You may be used to the notation 360° in place of 2π , but for the reasons explained in Appendix A and in footnote 15 of Ch. 3, this book tends to avoid the former notation.

substituting into the former yields

$$\psi_1 + \psi_2 + \psi_3 = \frac{2\pi}{2}, \quad (2.45)$$

which was to be demonstrated.

Extending the same technique to the case of an n -sided polygon, we have that

$$\begin{aligned} \sum_{k=1}^n \phi_k &= 2\pi, \\ \phi_k + \psi_k &= \frac{2\pi}{2}. \end{aligned}$$

Solving the latter equations for ϕ_k and substituting into the former yields

$$\sum_{k=1}^n \left(\frac{2\pi}{2} - \psi_k \right) = 2\pi,$$

or in other words

$$\sum_{k=1}^n \psi_k = (n-2) \frac{2\pi}{2}. \quad (2.46)$$

Equation (2.45) is then seen to be a special case of (2.46) with $n = 3$.

2.10 The Pythagorean theorem

Along with Euler's formula (5.11), the fundamental theorem of calculus (7.2) and Cauchy's integral formula (8.29), the Pythagorean theorem is one of the most famous results in all of mathematics. The theorem holds that

$$a^2 + b^2 = c^2, \quad (2.47)$$

where a , b and c are the lengths of the legs and diagonal of a right triangle, as in Fig. 2.3. Many proofs of the theorem are known.

One such proof posits a square of side length $a + b$ with a tilted square of side length c inscribed as in Fig. 2.4. The area of each of the four triangles in the figure is evidently $ab/2$. The area of the tilted inner square is c^2 . The area of the large outer square is $(a + b)^2$. But the large outer square is comprised of the tilted inner square plus the four triangles, hence the area

Figure 2.3: A right triangle.

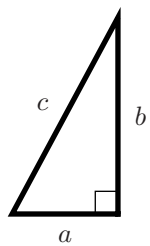
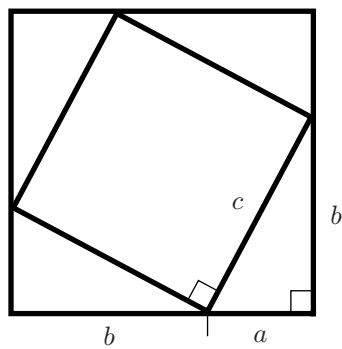


Figure 2.4: The Pythagorean theorem.



of the large outer square equals the area of the tilted inner square plus the areas of the four triangles. In mathematical symbols, this is

$$(a + b)^2 = c^2 + 4 \left(\frac{ab}{2} \right),$$

which simplifies directly to (2.47).²⁸

The Pythagorean theorem is readily extended to three dimensions as

$$a^2 + b^2 + h^2 = r^2, \quad (2.48)$$

where h is an altitude perpendicular to both a and b , thus also to c ; and where r is the corresponding three-dimensional diagonal: the diagonal of the right triangle whose legs are c and h . Inasmuch as (2.47) applies to any right triangle, it follows that $c^2 + h^2 = r^2$, which equation expands directly to yield (2.48).

2.11 Functions

This book is not the place for a gentle introduction to the concept of the function. Briefly, however, a *function* is a mapping from one number (or vector of several numbers) to another. For example, $f(x) = x^2 - 1$ is a function which maps 1 to 0 and -3 to 8, among others.

One often speaks of *domains* and *ranges* when discussing functions. The *domain* of a function is the set of numbers one can put into it. The *range* of a function is the corresponding set of numbers one can get out of it. In the example, if the domain is restricted to real x such that $|x| \leq 3$, then the corresponding range is $-1 \leq f(x) \leq 8$.

Other terms which arise when discussing functions are *root* (or *zero*), *singularity* and *pole*. A *root* (or *zero*) of a function is a domain point at which the function evaluates to zero (the example has roots at $x = \pm 1$). A *singularity* of a function is a domain point at which the function's output

²⁸This elegant proof is far simpler than the one famously given by the ancient geometer Euclid, yet more appealing than alternate proofs often found in print. Whether Euclid was acquainted with the simple proof given here this writer does not know, but it is possible [46, "Pythagorean theorem," 02:32, 31 March 2006] that Euclid chose his proof because it comported better with the restricted set of geometrical elements he permitted himself to work with. Be that as it may, the present writer encountered the proof this section gives somewhere years ago and has never seen it in print since, so can claim no credit for originating it. Unfortunately the citation is now long lost. A current source for the proof is [46] as cited earlier in this footnote.

diverges; that is, where the function's output is infinite.²⁹ A *pole* is a singularity that behaves locally like $1/x$ (rather than, for example, like $1/\sqrt{x}$). A singularity that behaves as $1/x^N$ is a *multiple pole*, which (§ 9.6.2) can be thought of as N poles. The example function $f(x)$ has no singularities for finite x ; however, the function $h(x) = 1/(x^2 - 1)$ has poles at $x = \pm 1$.

(Besides the root, the singularity and the pole, there is also the troublesome *branch point*, an infamous example of which is $z = 0$ in the function $g[z] = \sqrt{z}$. Branch points are important, but the book must lay a more extensive foundation before introducing them properly in § 8.5.³⁰)

2.12 Complex numbers (introduction)

Section 2.5.2 has introduced square roots. What it has not done is to tell us how to regard a quantity such as $\sqrt{-1}$. Since there exists no real number i such that

$$i^2 = -1 \tag{2.49}$$

and since the quantity i thus defined is found to be critically important across broad domains of higher mathematics, we accept (2.49) as the definition of a fundamentally new kind of quantity: *the imaginary number*.³¹

²⁹Here is one example of the book's deliberate lack of formal mathematical rigor. A more precise formalism to say that "the function's output is infinite" might be

$$\lim_{x \rightarrow x_0} |f(x)| = \infty.$$

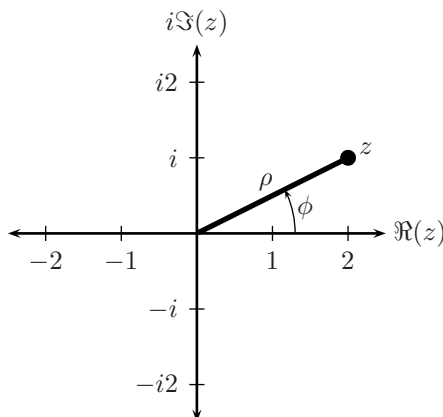
The applied mathematician tends to avoid such formalism where there seems no immediate use for it.

³⁰There is further the *essential singularity*, an example of which is $z = 0$ in $p(z) = \exp(1/z)$, but the best way to handle such unreasonable singularities is almost always to change a variable, as $w \leftarrow 1/z$, or otherwise to frame the problem such that one need not approach the singularity. This book will have little to say of such singularities.

³¹The English word *imaginary* is evocative, but perhaps not of quite the right concept in this usage. Imaginary numbers are not to mathematics as, say, imaginary elves are to the physical world. In the physical world, imaginary elves are (presumably) not substantial objects. However, in the mathematical realm, imaginary numbers *are* substantial. The word *imaginary* in the mathematical sense is thus more of a technical term than a descriptive adjective.

The number i is just a concept, of course, but then so is the number 1 (though you and I have often met one *of something*—one apple, one chair, one summer afternoon, etc.—neither of us has ever met just 1). The reason imaginary numbers are called "imaginary" probably has to do with the fact that they emerge from mathematical operations only, never directly from counting things. Notice, however, that the number $1/2$ never emerges directly from counting things, either. If for some reason the *i*year were offered as a unit of time, then the period separating your fourteenth and twenty-first birthdays could have

Figure 2.5: The complex (or Argand) plane, and a complex number z therein.



Imaginary numbers are given their own number line, plotted at right angles to the familiar real number line as in Fig. 2.5. The sum of a real number x and an imaginary number iy is the *complex number*

$$z = x + iy.$$

The *conjugate* z^* of this complex number is defined to be

$$z^* = x - iy.$$

The *magnitude* (or *modulus*, or *absolute value*) $|z|$ of the complex number is defined to be the length ρ in Fig. 2.5, which per the Pythagorean theorem (§ 2.10) is such that

$$|z|^2 = x^2 + y^2. \quad (2.50)$$

The *phase* $\arg z$ of the complex number is defined to be the angle ϕ in the figure, which in terms of the trigonometric functions of § 3.1³² is such that

$$\tan(\arg z) = \frac{y}{x}. \quad (2.51)$$

been measured as $-i7$ *iy*ears. Madness? No, let us not call it that; let us call it a useful formalism, rather.

The unpersuaded reader is asked to suspend judgment a while. He will soon see the use.

³²This is a forward reference. If the equation doesn't make sense to you yet for this reason, skip it for now. The important point is that $\arg z$ is the angle ϕ in the figure.

Specifically to extract the real and imaginary parts of a complex number, the notations

$$\begin{aligned}\Re(z) &= x, \\ \Im(z) &= y,\end{aligned}\tag{2.52}$$

are conventionally recognized (although often the symbols $\Re(\cdot)$ and $\Im(\cdot)$ are written $\text{Re}(\cdot)$ and $\text{Im}(\cdot)$, particularly when printed by hand).

2.12.1 Multiplication and division of complex numbers in rectangular form

Several elementary properties of complex numbers are readily seen if the fact that $i^2 = -1$ is kept in mind, including that

$$z_1 z_2 = (x_1 x_2 - y_1 y_2) + i(y_1 x_2 + x_1 y_2),\tag{2.53}$$

$$\begin{aligned}\frac{z_1}{z_2} &= \frac{x_1 + iy_1}{x_2 + iy_2} = \left(\frac{x_2 - iy_2}{x_2 - iy_2} \right) \frac{x_1 + iy_1}{x_2 + iy_2} \\ &= \frac{(x_1 x_2 + y_1 y_2) + i(y_1 x_2 - x_1 y_2)}{x_2^2 + y_2^2}.\end{aligned}\tag{2.54}$$

It is a curious fact that

$$\frac{1}{i} = -i.\tag{2.55}$$

It is a useful fact that

$$z^* z = x^2 + y^2 = |z|^2\tag{2.56}$$

(the curious fact, eqn. 2.55, is useful, too).

2.12.2 Complex conjugation

An important property of complex numbers descends subtly from the fact that

$$i^2 = -1 = (-i)^2.$$

If one defined some number $j \equiv -i$, claiming that j not i were the true imaginary unit,³³ then one would find that

$$(-j)^2 = -1 = j^2,$$

and thus that all the basic properties of complex numbers in the j system held just as well as they did in the i system. The units i and j would differ indeed, but would perfectly mirror one another in every respect.

³³[11, § I:22-5]

That is the basic idea. To establish it symbolically needs a page or so of slightly abstract algebra as follows, the goal of which will be to show that $[f(z)]^* = f(z^*)$ for some unspecified function $f(z)$ with specified properties. To begin with, if

$$z = x + iy,$$

then

$$z^* = x - iy$$

by definition. Proposing that $(z^{k-1})^* = (z^*)^{k-1}$ (which may or may not be true but for the moment we assume it), we can write

$$\begin{aligned} z^{k-1} &= s_{k-1} + it_{k-1}, \\ (z^*)^{k-1} &= s_{k-1} - it_{k-1}, \end{aligned}$$

where s_{k-1} and t_{k-1} are symbols introduced to represent the real and imaginary parts of z^{k-1} . Multiplying the former equation by $z = x + iy$ and the latter by $z^* = x - iy$, we have that

$$\begin{aligned} z^k &= (xs_{k-1} - yt_{k-1}) + i(ys_{k-1} + xt_{k-1}), \\ (z^*)^k &= (xs_{k-1} - yt_{k-1}) - i(ys_{k-1} + xt_{k-1}). \end{aligned}$$

With the definitions $s_k \equiv xs_{k-1} - yt_{k-1}$ and $t_k \equiv ys_{k-1} + xt_{k-1}$, this is written more succinctly

$$\begin{aligned} z^k &= s_k + it_k, \\ (z^*)^k &= s_k - it_k. \end{aligned}$$

In other words, if $(z^{k-1})^* = (z^*)^{k-1}$, then it necessarily follows that $(z^k)^* = (z^*)^k$. Solving the definitions of s_k and t_k for s_{k-1} and t_{k-1} yields the reverse definitions $s_{k-1} = (xs_k + yt_k)/(x^2 + y^2)$ and $t_{k-1} = (-ys_k + xt_k)/(x^2 + y^2)$. Therefore, except when $z = x + iy$ happens to be null or infinite, the implication is reversible by reverse reasoning, so by mathematical induction³⁴ we have that

$$(z^k)^* = (z^*)^k \tag{2.57}$$

for all integral k . We have also from (2.53) that

$$(z_1 z_2)^* = z_1^* z_2^* \tag{2.58}$$

³⁴*Mathematical induction* is an elegant old technique for the construction of mathematical proofs. Section 8.1 elaborates on the technique and offers a more extensive example. Beyond the present book, a very good introduction to mathematical induction is found in [18].

for any complex z_1 and z_2 .

Consequences of (2.57) and (2.58) include that if

$$f(z) \equiv \sum_{k=-\infty}^{\infty} (a_k + ib_k)(z - z_o)^k, \quad (2.59)$$

$$f^*(z) \equiv \sum_{k=-\infty}^{\infty} (a_k - ib_k)(z - z_o^*)^k, \quad (2.60)$$

where a_k and b_k are real and imaginary parts of the coefficients peculiar to the function $f(\cdot)$, then

$$[f(z)]^* = f^*(z^*). \quad (2.61)$$

In the common case where all $b_k = 0$ and $z_o = x_o$ is a real number, then $f(\cdot)$ and $f^*(\cdot)$ are the same function, so (2.61) reduces to the desired form

$$[f(z)]^* = f(z^*), \quad (2.62)$$

which says that the effect of conjugating the function's input is merely to conjugate its output.

Equation (2.62) expresses a significant, general rule of complex numbers and complex variables which is better explained in words than in mathematical symbols. The rule is this: for most equations and systems of equations used to model physical systems, *one can produce an equally valid alternate model simply by simultaneously conjugating all the complex quantities present.*³⁵

2.12.3 Power series and analytic functions (preview)

Equation (2.59) is a general power series in $z - z_o$.³⁶ Such power series have broad application.³⁷ It happens in practice that most functions of interest

³⁵[18][39]

³⁶[21, § 10.8]

³⁷That is a pretty impressive-sounding statement: "Such power series have broad application." However, air, words and molecules also have "broad application"; merely stating the fact does not tell us much. In fact the general power series is a sort of one-size-fits-all mathematical latex glove, which can be stretched to fit around almost any function. The interesting part is not so much in the general form (2.59) of the series as it is in the specific choice of a_k and b_k , which this section does not discuss.

Observe that the Taylor series (which this section also does not discuss; see § 8.3) is a power series with $a_k = b_k = 0$ for $k < 0$.

in modeling physical phenomena can conveniently be constructed as power series (or sums of power series)³⁸ with suitable choices of a_k , b_k and z_o .

The property (2.61) applies to all such functions, with (2.62) also applying to those for which $b_k = 0$ and $z_o = x_o$. The property the two equations represent is called the *conjugation property*. Basically, it says that if one replaces all the i in some mathematical model with $-i$, then the resulting conjugate model is equally as valid as the original.³⁹

Such functions, whether $b_k = 0$ and $z_o = x_o$ or not, are *analytic functions* (§ 8.4). In the formal mathematical definition, a function is analytic which is infinitely differentiable (Ch. 4) in the immediate domain neighborhood of interest. However, for applications a fair working definition of the analytic function might be “a function expressible as a power series.” Chapter 8 elaborates. All power series are infinitely differentiable except at their poles.

There nevertheless exist one common group of functions which cannot be constructed as power series. These all have to do with the parts of complex numbers and have been introduced in this very section: the magnitude $|\cdot|$; the phase $\arg(\cdot)$; the conjugate $(\cdot)^*$; and the real and imaginary parts $\Re(\cdot)$ and $\Im(\cdot)$. These functions are not analytic and do not in general obey the conjugation property. Also not analytic are the Heaviside unit step $u(t)$ and the Dirac delta $\delta(t)$ (§ 7.7), used to model discontinuities explicitly.

We shall have more to say about analytic functions in Ch. 8. We shall have more to say about complex numbers in §§ 3.11, 4.3.3, and 4.4, and much more yet in Ch. 5.

³⁸The careful reader might observe that this statement neglects the *Gibbs' phenomenon*, but that curious matter will be dealt with in [section not yet written].

³⁹To illustrate, from the fact that $(1 + i2)(2 + i3) + (1 - i) = -3 + i6$ the conjugation property infers immediately that $(1 - i2)(2 - i3) + (1 + i) = -3 - i6$. Observe however that no such property holds for the real parts: $(-1 + i2)(-2 + i3) + (-1 - i) \neq 3 + i6$.

Chapter 3

Trigonometry

Trigonometry is the branch of mathematics which relates angles to lengths. This chapter introduces the trigonometric functions and derives their several properties.

3.1 Definitions

Consider the circle-inscribed right triangle of Fig. 3.1.

In considering the circle, we shall find some terminology useful: the *angle* ϕ in the diagram is measured in *radians*, where a radian is the angle which, when centered in a *unit circle*, describes an arc of unit length.¹ Measured in radians, an angle ϕ intercepts an arc of curved length $\rho\phi$ on a circle of *radius* ρ (that is, of distance ρ from the circle's center to its perimeter). An angle in radians is a dimensionless number, so one need not write “ $\phi = 2\pi/4$ radians”; it suffices to write “ $\phi = 2\pi/4$.” In mathematical theory, we express angles in radians.

The angle of full revolution is given the symbol 2π —which thus is the circumference of a unit circle.² A quarter revolution, $2\pi/4$, is then the *right angle*, or square angle.

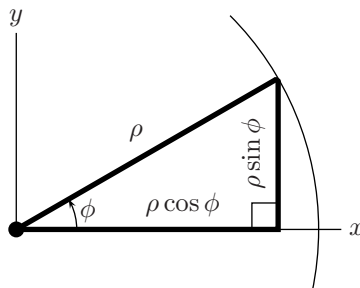
The trigonometric functions $\sin \phi$ and $\cos \phi$ (the “sine” and “cosine” of ϕ) relate the angle ϕ to the lengths shown in Fig. 3.1. The tangent function is then defined as

$$\tan \phi \equiv \frac{\sin \phi}{\cos \phi}, \quad (3.1)$$

¹The word “unit” means “one” in this context. A unit length is a length of 1 (not one centimeter or one mile, just an abstract 1). A unit circle is a circle of radius 1.

²Section 8.11 computes the numerical value of 2π .

Figure 3.1: The sine and the cosine (shown on a circle-inscribed right triangle, with the circle centered at the triangle's point).



which is the “rise” per unit “run,” or *slope*, of the triangle’s diagonal. Inverses of the three trigonometric functions can also be defined:

$$\begin{aligned}\arcsin(\sin \phi) &= \phi, \\ \arccos(\cos \phi) &= \phi, \\ \arctan(\tan \phi) &= \phi.\end{aligned}$$

When the last of these is written in the form

$$\arctan\left(\frac{y}{x}\right),$$

it is normally implied that x and y are to be interpreted as rectangular coordinates³ and that the arctan function is to return ϕ in the correct quadrant $-\pi < \phi \leq \pi$ (for example, $\arctan[1/(-1)] = [+3/8][2\pi]$, whereas $\arctan[(-1)/1] = [-1/8][2\pi]$). This is similarly the usual interpretation when an equation like

$$\tan \phi = \frac{y}{x}$$

is written.

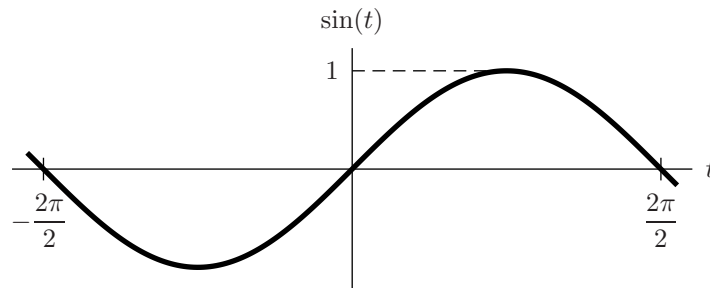
By the Pythagorean theorem (§ 2.10), it is seen generally that⁴

$$\cos^2 \phi + \sin^2 \phi = 1. \quad (3.2)$$

³*Rectangular coordinates* are pairs of numbers (x, y) which uniquely specify points in a plane. Conventionally, the x coordinate indicates distance eastward; the y coordinate, northward. For instance, the coordinates $(3, -4)$ mean the point three units eastward and four units southward (that is, -4 units northward) from the *origin* $(0, 0)$. A third rectangular coordinate can also be added— (x, y, z) —where the z indicates distance upward.

⁴The notation $\cos^2 \phi$ means $(\cos \phi)^2$.

Figure 3.2: The sine function.

Fig. 3.2 plots the sine function. The shape in the plot is called a *sinusoid*.

3.2 Simple properties

Inspecting Fig. 3.1 and observing (3.1) and (3.2), one readily discovers the several simple trigonometric properties of Table 3.1.

3.3 Scalars, vectors, and vector notation

In applied mathematics, a *vector* is a magnitude of some kind coupled with a direction.⁵ For example, “55 miles per hour northwestward” is a vector, as is the entity \mathbf{u} depicted in Fig. 3.3. The entity \mathbf{v} depicted in Fig. 3.4 is also a vector, in this case a three-dimensional one.

Many readers will already find the basic vector concept familiar, but for those who do not, a brief review: Vectors such as the

$$\begin{aligned}\mathbf{u} &= \hat{\mathbf{x}}x + \hat{\mathbf{y}}y, \\ \mathbf{v} &= \hat{\mathbf{x}}x + \hat{\mathbf{y}}y + \hat{\mathbf{z}}z\end{aligned}$$

⁵The same word *vector* is also used to indicate an ordered set of N scalars (§ 8.16) or an $N \times 1$ matrix (Ch. 11), but those are not the uses of the word meant here.

Table 3.1: Simple properties of the trigonometric functions.

$$\begin{array}{ll}
\sin(-\phi) &= -\sin \phi & \cos(-\phi) &= +\cos \phi \\
\sin(2\pi/4 - \phi) &= +\cos \phi & \cos(2\pi/4 - \phi) &= +\sin \phi \\
\sin(2\pi/2 - \phi) &= +\sin \phi & \cos(2\pi/2 - \phi) &= -\cos \phi \\
\sin(\phi \pm 2\pi/4) &= \pm \cos \phi & \cos(\phi \pm 2\pi/4) &= \mp \sin \phi \\
\sin(\phi \pm 2\pi/2) &= -\sin \phi & \cos(\phi \pm 2\pi/2) &= -\cos \phi \\
\sin(\phi + n2\pi) &= \sin \phi & \cos(\phi + n2\pi) &= \cos \phi
\end{array}$$

$$\begin{array}{ll}
\tan(-\phi) &= -\tan \phi \\
\tan(2\pi/4 - \phi) &= +1/\tan \phi \\
\tan(2\pi/2 - \phi) &= -\tan \phi \\
\tan(\phi \pm 2\pi/4) &= -1/\tan \phi \\
\tan(\phi \pm 2\pi/2) &= +\tan \phi \\
\tan(\phi + n2\pi) &= \tan \phi
\end{array}$$

$$\begin{array}{ll}
\frac{\sin \phi}{\cos \phi} &= \tan \phi \\
\cos^2 \phi + \sin^2 \phi &= 1 \\
1 + \tan^2 \phi &= \frac{1}{\cos^2 \phi} \\
1 + \frac{1}{\tan^2 \phi} &= \frac{1}{\sin^2 \phi}
\end{array}$$

Figure 3.3: A two-dimensional vector $\mathbf{u} = \hat{\mathbf{x}}x + \hat{\mathbf{y}}y$, shown with its rectangular components.

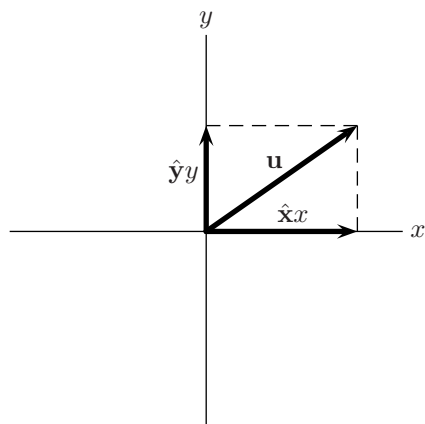
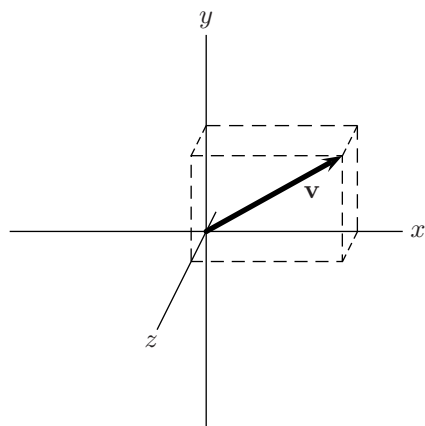


Figure 3.4: A three-dimensional vector $\mathbf{v} = \hat{\mathbf{x}}x + \hat{\mathbf{y}}y + \hat{\mathbf{z}}z$.



of the figures are composed of multiples of the *unit basis vectors* $\hat{\mathbf{x}}$, $\hat{\mathbf{y}}$ and $\hat{\mathbf{z}}$, which themselves are vectors of unit length pointing in the cardinal directions their respective symbols suggest.⁶ Any vector \mathbf{a} can be factored into a *magnitude* a and a *unit vector* $\hat{\mathbf{a}}$, as

$$\mathbf{a} = \hat{\mathbf{a}}a,$$

where the $\hat{\mathbf{a}}$ represents direction only and has unit magnitude by definition, and where the a represents magnitude only and carries the physical units if any.⁷ For example, $a = 55$ miles per hour, $\hat{\mathbf{a}} =$ northwestward. The unit vector $\hat{\mathbf{a}}$ itself can be expressed in terms of the unit basis vectors: for example, if $\hat{\mathbf{x}}$ points east and $\hat{\mathbf{y}}$ points north, then $\hat{\mathbf{a}} = -\hat{\mathbf{x}}(1/\sqrt{2}) + \hat{\mathbf{y}}(1/\sqrt{2})$, where per the Pythagorean theorem $(-1/\sqrt{2})^2 + (1/\sqrt{2})^2 = 1^2$.

A single number which is not a vector or a matrix (Ch. 11) is called a *scalar*. In the example, $a = 55$ miles per hour is a scalar. Though the scalar a in the example happens to be real, scalars can be complex, too—which might surprise one, since scalars by definition lack direction and the Argand phase ϕ of Fig. 2.5 so strongly resembles a direction. However, phase is not an actual direction in the vector sense (the real number line in the Argand plane cannot be said to run west-to-east, or anything like that). The x , y and z of Fig. 3.4 are each (possibly complex) scalars; $\mathbf{v} =$

⁶Printing by hand, one customarily writes a general vector like \mathbf{u} as “ \vec{u} ” or just “ \overline{u} ”, and a unit vector like $\hat{\mathbf{x}}$ as “ \hat{x} ”.

⁷The word “unit” here is unfortunately overloaded. As an adjective in mathematics, or in its noun form “unity,” it refers to the number one (1)—not one mile per hour, one kilogram, one Japanese yen or anything like that; just an abstract 1. The word “unit” itself as a noun however usually signifies a physical or financial reference quantity of measure, like a mile per hour, a kilogram or even a Japanese yen. There is no inherent mathematical unity to 1 mile per hour (otherwise known as 0.447 meters per second, among other names). By contrast, a “unitless 1”—a 1 with no physical unit attached—does represent mathematical unity.

Consider the ratio $r = h_1/h_o$ of your height h_1 to my height h_o . Maybe you are taller than I am and so $r = 1.05$ (not 1.05 cm or 1.05 feet, just 1.05). Now consider the ratio h_1/h_1 of your height to your own height. That ratio is of course unity, exactly 1.

There is nothing ephemeral in the concept of mathematical unity, nor in the concept of unitless quantities in general. The concept is quite straightforward and entirely practical. That $r > 1$ means neither more nor less than that you are taller than I am. In applications, one often puts physical quantities in ratio precisely to strip the physical units from them, comparing the ratio to unity without regard to physical units.

$\hat{\mathbf{x}}x + \hat{\mathbf{y}}y + \hat{\mathbf{z}}z$ is a vector. If x , y and z are complex, then⁸

$$\begin{aligned} |\mathbf{v}|^2 &= |x|^2 + |y|^2 + |z|^2 = x^*x + y^*y + z^*z \\ &= [\Re(x)]^2 + [\Im(x)]^2 + [\Re(y)]^2 + [\Im(y)]^2 \\ &\quad + [\Re(z)]^2 + [\Im(z)]^2. \end{aligned} \tag{3.3}$$

A point is sometimes identified by the vector expressing its distance and direction from the origin of the coordinate system. That is, the point (x, y) can be identified with the vector $\hat{\mathbf{x}}x + \hat{\mathbf{y}}y$. However, in the general case vectors are not associated with any particular origin; they represent distances and directions, not fixed positions.

Notice the relative orientation of the axes in Fig. 3.4. The axes are oriented such that if you point your flat right hand in the x direction, then bend your fingers in the y direction and extend your thumb, the thumb then points in the z direction. This is orientation by the *right-hand rule*. A left-handed orientation is equally possible, of course, but as neither orientation has a natural advantage over the other, we arbitrarily but conventionally accept the right-handed one as standard.⁹

3.4 Rotation

A fundamental problem in trigonometry arises when a vector

$$\mathbf{u} = \hat{\mathbf{x}}x + \hat{\mathbf{y}}y \tag{3.4}$$

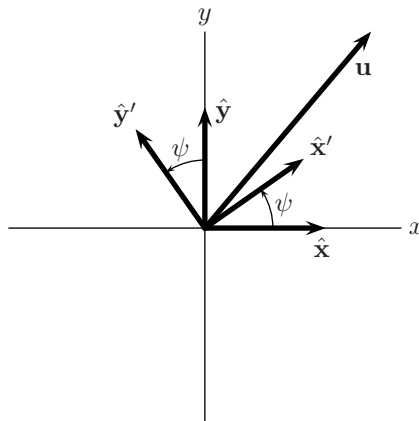
must be expressed in terms of alternate unit vectors $\hat{\mathbf{x}}'$ and $\hat{\mathbf{y}}'$, where $\hat{\mathbf{x}}'$ and $\hat{\mathbf{y}}'$ stand at right angles to one another and lie in the plane¹⁰ of $\hat{\mathbf{x}}$ and $\hat{\mathbf{y}}$,

⁸Some books print $|\mathbf{v}|$ as $\|\mathbf{v}\|$ or even $\|\mathbf{v}\|_2$ to emphasize that it represents the real, scalar magnitude of a complex vector. The reason the last notation subscripts a numeral 2 is obscure, having to do with the professional mathematician's generalized definition of a thing he calls the "norm." This book just renders it $|\mathbf{v}|$.

⁹The writer does not know the etymology for certain, but verbal lore in American engineering has it that the name "right-handed" comes from experience with a standard right-handed wood screw or machine screw. If you hold the screwdriver in your right hand and turn the screw in the natural manner clockwise, turning the screw slot from the x orientation toward the y , the screw advances away from you in the z direction into the wood or hole. If somehow you came across a left-handed screw, you'd probably find it easier to drive that screw with the screwdriver in your left hand.

¹⁰A *plane*, as the reader on this tier undoubtedly knows, is a flat (but not necessarily level) surface, infinite in extent unless otherwise specified. Space is three-dimensional. A plane is two-dimensional. A line is one-dimensional. A point is zero-dimensional. The plane belongs to this geometrical hierarchy.

Figure 3.5: Vector basis rotation.



but are rotated from the latter by an angle ψ as depicted in Fig. 3.5.¹¹ In terms of the trigonometric functions of § 3.1, evidently

$$\begin{aligned}\hat{\mathbf{x}}' &= +\hat{\mathbf{x}} \cos \psi + \hat{\mathbf{y}} \sin \psi, \\ \hat{\mathbf{y}}' &= -\hat{\mathbf{x}} \sin \psi + \hat{\mathbf{y}} \cos \psi;\end{aligned}\tag{3.5}$$

and by appeal to symmetry it stands to reason that

$$\begin{aligned}\hat{\mathbf{x}} &= +\hat{\mathbf{x}}' \cos \psi - \hat{\mathbf{y}}' \sin \psi, \\ \hat{\mathbf{y}} &= +\hat{\mathbf{x}}' \sin \psi + \hat{\mathbf{y}}' \cos \psi.\end{aligned}\tag{3.6}$$

Substituting (3.6) into (3.4) yields

$$\mathbf{u} = \hat{\mathbf{x}}'(x \cos \psi + y \sin \psi) + \hat{\mathbf{y}}'(-x \sin \psi + y \cos \psi),\tag{3.7}$$

which was to be derived.

Equation (3.7) finds general application where rotations in rectangular coordinates are involved. If the question is asked, “what happens if I rotate not the unit basis vectors but rather the vector \mathbf{u} instead?” the answer is

¹¹The “’” mark is pronounced “prime” or “primed” (for no especially good reason of which the author is aware, but anyway, that’s how it’s pronounced). Mathematical writing employs the mark for a variety of purposes. Here, the mark merely distinguishes the new unit vector $\hat{\mathbf{x}}'$ from the old $\hat{\mathbf{x}}$.

that it amounts to the same thing, except that the sense of the rotation is reversed:

$$\mathbf{u}' = \hat{\mathbf{x}}(x \cos \psi - y \sin \psi) + \hat{\mathbf{y}}(x \sin \psi + y \cos \psi). \quad (3.8)$$

Whether it is the basis or the vector which rotates thus depends on your point of view.¹²

3.5 Trigonometric functions of sums and differences of angles

With the results of § 3.4 in hand, we now stand in a position to consider trigonometric functions of sums and differences of angles. Let

$$\begin{aligned} \hat{\mathbf{a}} &\equiv \hat{\mathbf{x}} \cos \alpha + \hat{\mathbf{y}} \sin \alpha, \\ \hat{\mathbf{b}} &\equiv \hat{\mathbf{x}} \cos \beta + \hat{\mathbf{y}} \sin \beta, \end{aligned}$$

be vectors of unit length in the xy plane, respectively at angles α and β from the x axis. If we wanted $\hat{\mathbf{b}}$ to coincide with $\hat{\mathbf{a}}$, we would have to rotate it by $\psi = \alpha - \beta$. According to (3.8) and the definition of $\hat{\mathbf{b}}$, if we did this we would obtain

$$\begin{aligned} \hat{\mathbf{b}}' &= \hat{\mathbf{x}}[\cos \beta \cos(\alpha - \beta) - \sin \beta \sin(\alpha - \beta)] \\ &\quad + \hat{\mathbf{y}}[\cos \beta \sin(\alpha - \beta) + \sin \beta \cos(\alpha - \beta)]. \end{aligned}$$

Since we have deliberately chosen the angle of rotation such that $\hat{\mathbf{b}}' = \hat{\mathbf{a}}$, we can separately equate the $\hat{\mathbf{x}}$ and $\hat{\mathbf{y}}$ terms in the expressions for $\hat{\mathbf{a}}$ and $\hat{\mathbf{b}}'$ to obtain the pair of equations

$$\begin{aligned} \cos \alpha &= \cos \beta \cos(\alpha - \beta) - \sin \beta \sin(\alpha - \beta), \\ \sin \alpha &= \cos \beta \sin(\alpha - \beta) + \sin \beta \cos(\alpha - \beta). \end{aligned}$$

¹²This is only true, of course, with respect to the vectors themselves. When one actually rotates a physical body, the body experiences forces during rotation which might or might not change the body internally in some relevant way.

Solving the last pair simultaneously¹³ for $\sin(\alpha - \beta)$ and $\cos(\alpha - \beta)$ and observing that $\sin^2(\cdot) + \cos^2(\cdot) = 1$ yields

$$\begin{aligned}\sin(\alpha - \beta) &= \sin \alpha \cos \beta - \cos \alpha \sin \beta, \\ \cos(\alpha - \beta) &= \cos \alpha \cos \beta + \sin \alpha \sin \beta.\end{aligned}\tag{3.9}$$

With the change of variable $\beta \leftarrow -\beta$ and the observations from Table 3.1 that $\sin(-\phi) = -\sin \phi$ and $\cos(-\phi) = +\cos(\phi)$, eqns. (3.9) become

$$\begin{aligned}\sin(\alpha + \beta) &= \sin \alpha \cos \beta + \cos \alpha \sin \beta, \\ \cos(\alpha + \beta) &= \cos \alpha \cos \beta - \sin \alpha \sin \beta.\end{aligned}\tag{3.10}$$

Equations (3.9) and (3.10) are the basic formulas for trigonometric functions of sums and differences of angles.

3.5.1 Variations on the sums and differences

Several useful variations on (3.9) and (3.10) are achieved by combining the equations in various straightforward ways.¹⁴ These include

$$\begin{aligned}\sin \alpha \sin \beta &= \frac{\cos(\alpha - \beta) - \cos(\alpha + \beta)}{2}, \\ \sin \alpha \cos \beta &= \frac{\sin(\alpha - \beta) + \sin(\alpha + \beta)}{2}, \\ \cos \alpha \cos \beta &= \frac{\cos(\alpha - \beta) + \cos(\alpha + \beta)}{2}.\end{aligned}\tag{3.11}$$

¹³The easy way to do this is

- to subtract $\sin \beta$ times the first equation from $\cos \beta$ times the second, then to solve the result for $\sin(\alpha - \beta)$;
- to add $\cos \beta$ times the first equation to $\sin \beta$ times the second, then to solve the result for $\cos(\alpha - \beta)$.

This shortcut technique for solving a pair of equations simultaneously for a pair of variables is well worth mastering. In this book alone, it proves useful many times.

¹⁴Refer to footnote 13 above for the technique.

With the change of variables $\delta \leftarrow \alpha - \beta$ and $\gamma \leftarrow \alpha + \beta$, (3.9) and (3.10) become

$$\begin{aligned}\sin \delta &= \sin \left(\frac{\gamma + \delta}{2} \right) \cos \left(\frac{\gamma - \delta}{2} \right) - \cos \left(\frac{\gamma + \delta}{2} \right) \sin \left(\frac{\gamma - \delta}{2} \right), \\ \cos \delta &= \cos \left(\frac{\gamma + \delta}{2} \right) \cos \left(\frac{\gamma - \delta}{2} \right) + \sin \left(\frac{\gamma + \delta}{2} \right) \sin \left(\frac{\gamma - \delta}{2} \right), \\ \sin \gamma &= \sin \left(\frac{\gamma + \delta}{2} \right) \cos \left(\frac{\gamma - \delta}{2} \right) + \cos \left(\frac{\gamma + \delta}{2} \right) \sin \left(\frac{\gamma - \delta}{2} \right), \\ \cos \gamma &= \cos \left(\frac{\gamma + \delta}{2} \right) \cos \left(\frac{\gamma - \delta}{2} \right) - \sin \left(\frac{\gamma + \delta}{2} \right) \sin \left(\frac{\gamma - \delta}{2} \right).\end{aligned}$$

Combining these in various ways, we have that

$$\begin{aligned}\sin \gamma + \sin \delta &= 2 \sin \left(\frac{\gamma + \delta}{2} \right) \cos \left(\frac{\gamma - \delta}{2} \right), \\ \sin \gamma - \sin \delta &= 2 \cos \left(\frac{\gamma + \delta}{2} \right) \sin \left(\frac{\gamma - \delta}{2} \right), \\ \cos \delta + \cos \gamma &= 2 \cos \left(\frac{\gamma + \delta}{2} \right) \cos \left(\frac{\gamma - \delta}{2} \right), \\ \cos \delta - \cos \gamma &= 2 \sin \left(\frac{\gamma + \delta}{2} \right) \sin \left(\frac{\gamma - \delta}{2} \right).\end{aligned}\tag{3.12}$$

3.5.2 Trigonometric functions of double and half angles

If $\alpha = \beta$, then eqns. (3.10) become the *double-angle formulas*

$$\begin{aligned}\sin 2\alpha &= 2 \sin \alpha \cos \alpha, \\ \cos 2\alpha &= 2 \cos^2 \alpha - 1 = \cos^2 \alpha - \sin^2 \alpha = 1 - 2 \sin^2 \alpha.\end{aligned}\tag{3.13}$$

Solving (3.13) for $\sin^2 \alpha$ and $\cos^2 \alpha$ yields the *half-angle formulas*

$$\begin{aligned}\sin^2 \alpha &= \frac{1 - \cos 2\alpha}{2}, \\ \cos^2 \alpha &= \frac{1 + \cos 2\alpha}{2}.\end{aligned}\tag{3.14}$$

3.6 Trigonometric functions of the hour angles

In general one uses the Taylor series of Ch. 8 to calculate trigonometric functions of specific angles. However, for angles which happen to be integral

multiples of an *hour*—there are twenty-four or 0x18 hours in a circle, just as there are twenty-four or 0x18 hours in a day¹⁵—for such angles simpler expressions exist. Figure 3.6 shows the angles. Since such angles arise very frequently in practice, it seems worth our while to study them specially.

Table 3.2 tabulates the trigonometric functions of these *hour angles*. To see how the values in the table are calculated, look at the square and the equilateral triangle¹⁶ of Fig. 3.7. Each of the square’s four angles naturally measures six hours; and since a triangle’s angles always total twelve hours (§ 2.9.3), by symmetry each of the angles of the equilateral triangle in the figure measures four. Also by symmetry, the perpendicular splits the triangle’s top angle into equal halves of two hours each and its bottom leg into equal segments of length 1/2 each; and the diagonal splits the square’s corner into equal halves of three hours each. The Pythagorean theorem (§ 2.10) then supplies the various other lengths in the figure, after which we observe from Fig. 3.1 that

- the sine of a non-right angle in a right triangle is the opposite leg’s length divided by the diagonal’s,

¹⁵Hence an hour is 15° , but you weren’t going to write your angles in such inelegant conventional notation as “ 15° ,” were you? Well, if you were, you’re in good company.

The author is fully aware of the barrier the unfamiliar notation poses for most first-time readers of the book. The barrier is erected neither lightly nor disrespectfully. Consider:

- There are 0x18 hours in a circle.
- There are 360 degrees in a circle.

Both sentences say the same thing, don’t they? But even though the “0x” hex prefix is a bit clumsy, the first sentence nevertheless says the thing rather better. The reader is urged to invest the attention and effort to master the notation.

There is a psychological trap regarding the hour. The familiar, standard clock face shows only twelve hours not twenty-four, so the angle between eleven o’clock and twelve *on the clock face* is not an hour of arc! That angle is two hours of arc. This is so because the clock face’s geometry is artificial. If you have ever been to the Old Royal Observatory at Greenwich, England, you may have seen the big clock face there with all twenty-four hours on it. It’d be a bit hard to read the time from such a crowded clock face were it not so big, but anyway, the angle between hours on the Greenwich clock is indeed an honest hour of arc. [3]

The hex and hour notations are recommended mostly only for theoretical math work. It is not claimed that they offer much benefit in most technical work of the less theoretical kinds. If you wrote an engineering memo describing a survey angle as 0x1.80 hours instead of 22.5 degrees, for example, you’d probably not like the reception the memo got. Nonetheless, the improved notation fits a book of this kind so well that the author hazards it. It is hoped that after trying the notation a while, the reader will approve the choice.

¹⁶An *equilateral* triangle is, as the name and the figure suggest, a triangle whose three sides all have the same length.

Figure 3.6: The 0x18 hours in a circle.

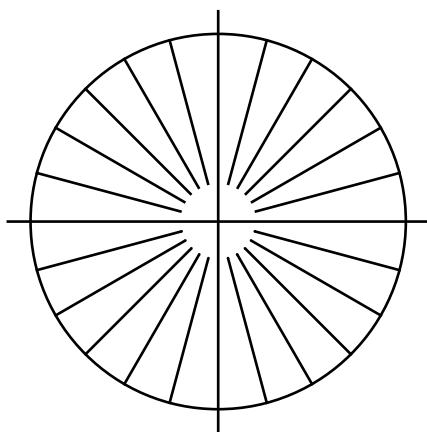


Figure 3.7: A square and an equilateral triangle for calculating trigonometric functions of the hour angles.

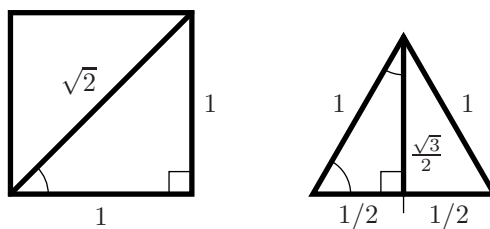
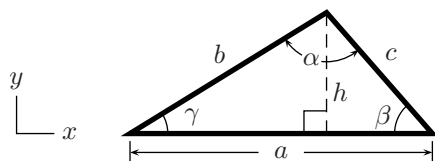


Table 3.2: Trigonometric functions of the hour angles.

ANGLE ϕ		$\sin \phi$	$\tan \phi$	$\cos \phi$
[radians]	[hours]			
0	0	0	0	1
$\frac{2\pi}{0x18}$	1	$\frac{\sqrt{3}-1}{2\sqrt{2}}$	$\frac{\sqrt{3}-1}{\sqrt{3}+1}$	$\frac{\sqrt{3}+1}{2\sqrt{2}}$
$\frac{2\pi}{0xC}$	2	$\frac{1}{2}$	$\frac{1}{\sqrt{3}}$	$\frac{\sqrt{3}}{2}$
$\frac{2\pi}{8}$	3	$\frac{1}{\sqrt{2}}$	1	$\frac{1}{\sqrt{2}}$
$\frac{2\pi}{6}$	4	$\frac{\sqrt{3}}{2}$	$\sqrt{3}$	$\frac{1}{2}$
$\frac{(5)(2\pi)}{0x18}$	5	$\frac{\sqrt{3}+1}{2\sqrt{2}}$	$\frac{\sqrt{3}+1}{\sqrt{3}-1}$	$\frac{\sqrt{3}-1}{2\sqrt{2}}$
$\frac{2\pi}{4}$	6	1	∞	0

Figure 3.8: The laws of sines and cosines.



- the tangent is the opposite leg's length divided by the adjacent leg's, and
- the cosine is the adjacent leg's length divided by the diagonal's.

With this observation and the lengths in the figure, one can calculate the sine, tangent and cosine of angles of two, three and four hours.

The values for one and five hours are found by applying (3.9) and (3.10) against the values for two and three hours just calculated. The values for zero and six hours are, of course, seen by inspection.¹⁷

3.7 The laws of sines and cosines

Refer to the triangle of Fig. 3.8. By the definition of the sine function, one can write that

$$c \sin \beta = h = b \sin \gamma,$$

or in other words that

$$\frac{\sin \beta}{b} = \frac{\sin \gamma}{c}.$$

¹⁷The creative reader may notice that he can extend the table to any angle by repeated application of the various sum, difference and half-angle formulas from the preceding sections to the values already in the table. However, the Taylor series (§ 8.9) offers a cleaner, quicker way to calculate trigonometrics of non-hour angles.

But there is nothing special about β and γ ; what is true for them must be true for α , too.¹⁸ Hence,

$$\frac{\sin \alpha}{a} = \frac{\sin \beta}{b} = \frac{\sin \gamma}{c}. \quad (3.15)$$

This equation is known as *the law of sines*.

On the other hand, if one expresses a and b as vectors emanating from the point γ ,¹⁹

$$\begin{aligned} \mathbf{a} &= \hat{\mathbf{x}}a, \\ \mathbf{b} &= \hat{\mathbf{x}}b \cos \gamma + \hat{\mathbf{y}}b \sin \gamma, \end{aligned}$$

then

$$\begin{aligned} c^2 &= |\mathbf{b} - \mathbf{a}|^2 \\ &= (b \cos \gamma - a)^2 + (b \sin \gamma)^2 \\ &= a^2 + (b^2)(\cos^2 \gamma + \sin^2 \gamma) - 2ab \cos \gamma. \end{aligned}$$

Since $\cos^2(\cdot) + \sin^2(\cdot) = 1$, this is

$$c^2 = a^2 + b^2 - 2ab \cos \gamma, \quad (3.16)$$

known as *the law of cosines*.

3.8 Summary of properties

Table 3.2 on page 58 has listed the values of trigonometric functions of the hour angles. Table 3.1 on page 48 has summarized simple properties of the trigonometric functions. Table 3.3 summarizes further properties, gathering them from §§ 3.4, 3.5 and 3.7.

¹⁸“But,” it is objected, “there *is* something special about α . The perpendicular h drops from it.”

True. However, the h is just a utility variable to help us to manipulate the equation into the desired form; we’re not interested in h itself. Nothing prevents us from dropping additional perpendiculars h_β and h_γ from the other two corners and using those as utility variables, too, if we like. We can use any utility variables we want.

¹⁹Here is another example of the book’s judicious relaxation of formal rigor. Of course there is no “point γ ”; γ is an angle not a point. However, the writer suspects in light of Fig. 3.8 that few readers will be confused as to which point is meant. The skillful applied mathematician does not multiply labels without need.

Table 3.3: Further properties of the trigonometric functions.

$$\begin{aligned}
\mathbf{u} &= \hat{\mathbf{x}}'(x \cos \psi + y \sin \psi) + \hat{\mathbf{y}}'(-x \sin \psi + y \cos \psi) \\
\sin(\alpha \pm \beta) &= \sin \alpha \cos \beta \pm \cos \alpha \sin \beta \\
\cos(\alpha \pm \beta) &= \cos \alpha \cos \beta \mp \sin \alpha \sin \beta \\
\sin \alpha \sin \beta &= \frac{\cos(\alpha - \beta) - \cos(\alpha + \beta)}{2} \\
\sin \alpha \cos \beta &= \frac{\sin(\alpha - \beta) + \sin(\alpha + \beta)}{2} \\
\cos \alpha \cos \beta &= \frac{\cos(\alpha - \beta) + \cos(\alpha + \beta)}{2} \\
\sin \gamma + \sin \delta &= 2 \sin \left(\frac{\gamma + \delta}{2} \right) \cos \left(\frac{\gamma - \delta}{2} \right) \\
\sin \gamma - \sin \delta &= 2 \cos \left(\frac{\gamma + \delta}{2} \right) \sin \left(\frac{\gamma - \delta}{2} \right) \\
\cos \delta + \cos \gamma &= 2 \cos \left(\frac{\gamma + \delta}{2} \right) \cos \left(\frac{\gamma - \delta}{2} \right) \\
\cos \delta - \cos \gamma &= 2 \sin \left(\frac{\gamma + \delta}{2} \right) \sin \left(\frac{\gamma - \delta}{2} \right) \\
\sin 2\alpha &= 2 \sin \alpha \cos \alpha \\
\cos 2\alpha &= 2 \cos^2 \alpha - 1 = \cos^2 \alpha - \sin^2 \alpha = 1 - 2 \sin^2 \alpha \\
\sin^2 \alpha &= \frac{1 - \cos 2\alpha}{2} \\
\cos^2 \alpha &= \frac{1 + \cos 2\alpha}{2} \\
\frac{\sin \gamma}{c} &= \frac{\sin \alpha}{a} = \frac{\sin \beta}{b} \\
c^2 &= a^2 + b^2 - 2ab \cos \gamma
\end{aligned}$$

3.9 Cylindrical and spherical coordinates

Section 3.3 has introduced the concept of the vector

$$\mathbf{v} = \hat{\mathbf{x}}x + \hat{\mathbf{y}}y + \hat{\mathbf{z}}z.$$

The coefficients (x, y, z) on the equation's right side are *coordinates*—specifically, *rectangular coordinates*²⁰—which given a specific orthonormal set of unit basis vectors $[\hat{\mathbf{x}} \ \hat{\mathbf{y}} \ \hat{\mathbf{z}}]$ uniquely identify a point (see Fig. 3.4 on page 49). Such rectangular coordinates are simple and general, and are convenient for many purposes. However, there are at least two broad classes of conceptually simple problems for which rectangular coordinates tend to be inconvenient: problems in which an axis or a point dominates. Consider for example an electric wire's magnetic field, whose intensity varies with distance from the wire (an axis); or the illumination a lamp sheds on a printed page of this book, which depends on the book's distance from the lamp (a point).

To attack a problem dominated by an axis, the *cylindrical coordinates* $(\rho; \phi, z)$ can be used instead of the rectangular coordinates (x, y, z) . To attack a problem dominated by a point, the *spherical coordinates* $(r; \theta; \phi)$ can be used.²¹ Refer to Fig. 3.9. Such coordinates are related to one another and to the rectangular coordinates by the formulas of Table 3.4.

Cylindrical and spherical coordinates can greatly simplify the analyses of the kinds of problems they respectively fit, but they come at a price. There are no constant unit basis vectors to match them. That is,

$$\mathbf{v} = \hat{\mathbf{x}}x + \hat{\mathbf{y}}y + \hat{\mathbf{z}}z \neq \hat{\boldsymbol{\rho}}\rho + \hat{\boldsymbol{\phi}}\phi + \hat{\mathbf{z}}z \neq \hat{\mathbf{r}}r + \hat{\boldsymbol{\theta}}\theta + \hat{\boldsymbol{\phi}}\phi.$$

It doesn't work that way. Nevertheless *variable* unit basis vectors are defined:

$$\begin{aligned}\hat{\boldsymbol{\rho}} &\equiv +\hat{\mathbf{x}} \cos \phi + \hat{\mathbf{y}} \sin \phi, \\ \hat{\boldsymbol{\phi}} &\equiv -\hat{\mathbf{x}} \sin \phi + \hat{\mathbf{y}} \cos \phi, \\ \hat{\mathbf{r}} &\equiv +\hat{\mathbf{z}} \cos \theta + \hat{\boldsymbol{\rho}} \sin \theta, \\ \hat{\boldsymbol{\theta}} &\equiv -\hat{\mathbf{z}} \sin \theta + \hat{\boldsymbol{\rho}} \cos \theta.\end{aligned}\tag{3.17}$$

²⁰*Orthonormal* in this context means “of unit length and at right angles to the other vectors in the set.” [46, “Orthonormality,” 14:19, 7 May 2006]

²¹Notice that the ϕ is conventionally written second in cylindrical $(\rho; \phi, z)$ but third in spherical $(r; \theta; \phi)$ coordinates. This odd-seeming convention is to maintain proper right-handed coordinate rotation. (The explanation will seem clearer once Ch. [not yet written] is read.)

Figure 3.9: A point on a sphere, in spherical $(r; \theta; \phi)$ and cylindrical $(\rho; \phi, z)$ coordinates. (The axis labels bear circumflexes in this figure only to disambiguate the \hat{z} axis from the cylindrical coordinate z .)

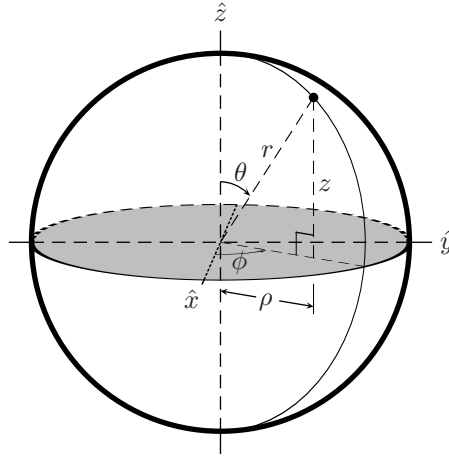


Table 3.4: Relations among the rectangular, cylindrical and spherical coordinates.

$$\begin{aligned}
 \rho^2 &= x^2 + y^2 \\
 r^2 &= \rho^2 + z^2 = x^2 + y^2 + z^2 \\
 \tan \theta &= \frac{\rho}{z} \\
 \tan \phi &= \frac{y}{x} \\
 z &= r \cos \theta \\
 \rho &= r \sin \theta \\
 x &= \rho \cos \phi = r \sin \theta \cos \phi \\
 y &= \rho \sin \phi = r \sin \theta \sin \phi
 \end{aligned}$$

Such variable unit basis vectors point locally in the directions in which their respective coordinates advance.

Convention usually orients $\hat{\mathbf{z}}$ in the direction of a problem's axis. Occasionally however a problem arises in which it is more convenient to orient $\hat{\mathbf{x}}$ or $\hat{\mathbf{y}}$ in the direction of the problem's axis (usually because $\hat{\mathbf{z}}$ has already been established in the direction of some other pertinent axis). Changing the meanings of known symbols like ρ , θ and ϕ is usually not a good idea, but you can use symbols like

$$\begin{aligned}\rho_x^2 &= y^2 + z^2, & \rho_y^2 &= z^2 + x^2, \\ \tan \theta_x &= \frac{\rho_x}{x}, & \tan \theta_y &= \frac{\rho_y}{y}, \\ \tan \phi_x &= \frac{z}{y}, & \tan \phi_y &= \frac{x}{z},\end{aligned}\tag{3.18}$$

instead if needed.²²

3.10 The complex triangle inequalities

If \mathbf{a} , \mathbf{b} and \mathbf{c} represent the three sides of a triangle such that $\mathbf{a} + \mathbf{b} + \mathbf{c} = 0$, then per (2.44)

$$|\mathbf{a}| - |\mathbf{b}| \leq |\mathbf{a} + \mathbf{b}| \leq |\mathbf{a}| + |\mathbf{b}|.$$

These are just the triangle inequalities of § 2.9.2 in vector notation.²³ But if the triangle inequalities hold for vectors in a plane, then why not equally for complex numbers? Consider the geometric interpretation of the Argand plane of Fig. 2.5 on page 40. Evidently,

$$|z_1| - |z_2| \leq |z_1 + z_2| \leq |z_1| + |z_2|\tag{3.19}$$

for any two complex numbers z_1 and z_2 . Extending the sum inequality, we have that

$$\left| \sum_k z_k \right| \leq \sum_k |z_k|.\tag{3.20}$$

Naturally, (3.19) and (3.20) hold equally well for real numbers as for complex; one may find the latter formula useful for sums of real numbers, for

²²Symbols like ρ_x are logical but, as far as this writer is aware, not standard. The writer is not aware of any conventionally established symbols for quantities like these.

²³Reading closely, one might note that § 2.9.2 uses the “<” sign rather than the “≤,” but that’s all right. See § 1.2.2.

example, when some of the numbers summed are positive and others negative.

An important consequence of (3.20) is that if $\sum |z_k|$ converges, then $\sum z_k$ also converges. Such a consequence is important because mathematical derivations sometimes need the convergence of $\sum z_k$ established, which can be hard to do directly. Convergence of $\sum |z_k|$, which per (3.20) implies convergence of $\sum z_k$, is often easier to establish.

Equation (3.20) will find use among other places in § 8.10.3.

3.11 De Moivre's theorem

Compare the Argand-plotted complex number of Fig. 2.5 (page 40) against the vector of Fig. 3.3 (page 49). Although complex numbers are scalars not vectors, the figures do suggest an analogy between complex phase and vector direction. With reference to Fig. 2.5 we can write

$$z = (\rho)(\cos \phi + i \sin \phi) = \rho \operatorname{cis} \phi, \quad (3.21)$$

where

$$\operatorname{cis} \phi \equiv \cos \phi + i \sin \phi. \quad (3.22)$$

If $z = x + iy$, then evidently

$$\begin{aligned} x &= \rho \cos \phi, \\ y &= \rho \sin \phi. \end{aligned} \quad (3.23)$$

Per (2.53),

$$z_1 z_2 = (x_1 x_2 - y_1 y_2) + i(y_1 x_2 + x_1 y_2).$$

Applying (3.23) to the equation yields

$$\frac{z_1 z_2}{\rho_1 \rho_2} = (\cos \phi_1 \cos \phi_2 - \sin \phi_1 \sin \phi_2) + i(\sin \phi_1 \cos \phi_2 + \cos \phi_1 \sin \phi_2).$$

But according to (3.10), this is just

$$\frac{z_1 z_2}{\rho_1 \rho_2} = \cos(\phi_1 + \phi_2) + i \sin(\phi_1 + \phi_2),$$

or in other words

$$z_1 z_2 = \rho_1 \rho_2 \operatorname{cis}(\phi_1 + \phi_2). \quad (3.24)$$

Equation (3.24) is an important result. It says that if you want to multiply complex numbers, it suffices

- to multiply their magnitudes and
- to add their phases.

It follows by parallel reasoning (or by extension) that

$$\frac{z_1}{z_2} = \frac{\rho_1}{\rho_2} \operatorname{cis}(\phi_1 - \phi_2) \quad (3.25)$$

and by extension that

$$z^a = \rho^a \operatorname{cis} a\phi. \quad (3.26)$$

Equations (3.24), (3.25) and (3.26) are known as *de Moivre's theorem*.^{24,25}

We have not shown yet, but shall in § 5.4, that

$$\operatorname{cis} \phi = \exp i\phi = e^{i\phi},$$

where $\exp(\cdot)$ is the natural exponential function and e is the natural logarithmic base, both defined in Ch. 5. De Moivre's theorem is most useful in light of this.

²⁴Also called *de Moivre's formula*. Some authors apply the name of de Moivre directly only to (3.26), or to some variation thereof; but since the three equations express essentially the same idea, if you refer to any of them as *de Moivre's theorem* then you are unlikely to be misunderstood.

²⁵[39][46]

Chapter 4

The derivative

The mathematics of *calculus* concerns a complementary pair of questions:¹

- Given some function $f(t)$, what is the function's instantaneous rate of change, or *derivative*, $f'(t)$?
- Interpreting some function $f'(t)$ as an instantaneous rate of change, what is the corresponding accretion, or *integral*, $f(t)$?

This chapter builds toward a basic understanding of the first question.

4.1 Infinitesimals and limits

Calculus systematically treats numbers so large and so small, they lie beyond the reach of our mundane number system.

¹Although once grasped the concept is relatively simple, to understand this pair of questions, so briefly stated, is no trivial thing. They are the pair which eluded or confounded the most brilliant mathematical minds of the ancient world.

The greatest conceptual hurdle—the stroke of brilliance—probably lies simply in stating the pair of questions clearly. Sir Isaac Newton and G.W. Leibnitz cleared this hurdle for us in the seventeenth century, so now at least we know the right pair of questions to ask. With the pair in hand, the calculus beginner's first task is quantitatively to understand the pair's interrelationship, generality and significance. Such an understanding constitutes the basic calculus concept.

It cannot be the role of a book like this one to lead the beginner gently toward an apprehension of the basic calculus concept. Once grasped, the concept is simple and briefly stated. In this book we necessarily state the concept briefly, then move along. Many instructional textbooks—[18] is a worthy example—have been written to lead the beginner gently. Although a sufficiently talented, dedicated beginner could perhaps obtain the basic calculus concept directly here, he would probably find it quicker and more pleasant to begin with a book like the one referenced.

4.1.1 The infinitesimal

A number ϵ is an *infinitesimal* if it is so small that

$$0 < |\epsilon| < a$$

for all possible mundane positive numbers a .

This is somewhat a difficult concept, so if it is not immediately clear then let us approach the matter colloquially. Let me propose to you that I have an infinitesimal.

“How big is your infinitesimal?” you ask.

“Very, very small,” I reply.

“How small?”

“Very small.”

“Smaller than 0x0.01?”

“Smaller than what?”

“Than 2^{-8} . You said that we should use hexadecimal notation in this book, remember?”

“Sorry. Yes, right, smaller than 0x0.01.”

“What about 0x0.0001? Is it smaller than that?”

“Much smaller.”

“Smaller than 0x0.0000 0000 0000 0001?”

“Smaller.”

“Smaller than $2^{-0x1\ 0000\ 0000\ 0000\ 0000}$?”

“Now *that* is an impressively small number. Nevertheless, my infinitesimal is smaller still.”

“Zero, then.”

“Oh, no. Bigger than that. My infinitesimal is definitely bigger than zero.”

This is the idea of the infinitesimal. It is a definite number of a certain nonzero magnitude, but its smallness conceptually lies beyond the reach of our mundane number system.

If ϵ is an infinitesimal, then $1/\epsilon$ can be regarded as an *infinity*: a very large number much larger than any mundane number one can name.

The principal advantage of using symbols like ϵ rather than 0 for infinitesimals is in that it permits us conveniently to compare one infinitesimal against another, to add them together, to divide them, etc. For instance, if $\delta = 3\epsilon$ is another infinitesimal, then the quotient δ/ϵ is not some unfathomable $0/0$; rather it is $\delta/\epsilon = 3$. In physical applications, the infinitesimals are often not true mathematical infinitesimals but rather relatively very small quantities such as the mass of a wood screw compared to the mass

of a wooden house frame, or the audio power of your voice compared to that of a jet engine. The additional cost of inviting one more guest to the wedding may or may not be infinitesimal, depending on your point of view. The key point is that the infinitesimal quantity be negligible by comparison, whatever “negligible” means in the context.²

The second-order infinitesimal ϵ^2 is so small on the scale of the common, first-order infinitesimal ϵ that the even latter cannot measure it. The ϵ^2 is an infinitesimal to the infinitesimals. Third- and higher-order infinitesimals are likewise possible.

The notation $u \ll v$, or $v \gg u$, indicates that u is much less than v , typically such that one can regard the quantity u/v to be an infinitesimal. In fact, one common way to specify that ϵ be infinitesimal is to write that $\epsilon \ll 1$.

4.1.2 Limits

The notation $\lim_{z \rightarrow z_o}$ indicates that z draws as near to z_o as it possibly can. When written $\lim_{z \rightarrow z_o^+}$, the implication is that z draws toward z_o from the positive side such that $z > z_o$. Similarly, when written $\lim_{z \rightarrow z_o^-}$, the implication is that z draws toward z_o from the negative side.

The reason for the notation is to provide a way to handle expressions like

$$\frac{3z}{2z}$$

as z vanishes:

$$\lim_{z \rightarrow 0} \frac{3z}{2z} = \frac{3}{2}.$$

The symbol “ \lim_Q ” is short for “in the limit as Q .”

Notice that \lim is not a function like \log or \sin . It is just a reminder that a quantity approaches some value, used when saying that the quantity

²Among scientists and engineers who study wave phenomena, there is an old rule of thumb that sinusoidal waveforms be discretized not less finely than ten points per wavelength. In keeping with this book’s adecimal theme (Appendix A) and the concept of the hour of arc (§ 3.6), we should probably render the rule as *twelve* points per wavelength here. In any case, even very roughly speaking, a quantity greater than $1/0xC$ of the principal to which it compares probably cannot rightly be regarded as infinitesimal. On the other hand, a quantity less than $1/0x10000$ of the principal is indeed infinitesimal for most practical purposes (but not all: for example, positions of spacecraft and concentrations of chemical impurities must sometimes be accounted more precisely). For quantities between $1/0xC$ and $1/0x10000$, it depends on the accuracy one seeks.

equaled the value would be confusing. Consider that to say

$$\lim_{z \rightarrow 2^-} (z + 2) = 4$$

is just a fancy way of saying that $2 + 2 = 4$. The \lim notation is convenient to use sometimes, but it is not magical. Don't let it confuse you.

4.2 Combinatorics

In its general form, the problem of selecting k specific items out of a set of n available items belongs to probability theory ([chapter not yet written]). In its basic form, however, the same problem also applies to the handling of polynomials or power series. This section treats the problem in its basic form.³

4.2.1 Combinations and permutations

Consider the following scenario. I have several small wooden blocks of various shapes and sizes, painted different colors so that you can clearly tell each block from the others. If I offer you the blocks and you are free to take all, some or none of them at your option, if you can take whichever blocks you want, then how many distinct choices of blocks do you have? Answer: you have 2^n choices, because you can accept or reject the first block, then accept or reject the second, then the third, and so on.

Now, suppose that what you want is exactly k blocks, neither more nor fewer. Desiring exactly k blocks, you select your favorite block first: there are n options for this. Then you select your second favorite: for this, there are $n - 1$ options (why not n options? because you have already taken one block from me; I have only $n - 1$ blocks left). Then you select your third favorite—for this there are $n - 2$ options—and so on until you have k blocks. There are evidently

$$P \binom{n}{k} \equiv n!/(n - k)! \quad (4.1)$$

ordered ways, or *permutations*, available for you to select exactly k blocks.

However, some of these distinct permutations put exactly the same *combination* of blocks in your hand; for instance, the permutations red-green-blue and green-red-blue constitute the same combination, whereas red-white-blue is a different combination entirely. For a single combination

³[18]

of k blocks (red, green, blue), evidently $k!$ permutations are possible (red-green-blue, red-blue-green, green-red-blue, green-blue-red, blue-red-green, blue-green-red). Hence dividing the number of permutations (4.1) by $k!$ yields the number of combinations

$$\binom{n}{k} \equiv \frac{n!/(n-k)!}{k!}. \quad (4.2)$$

Properties of the number $\binom{n}{k}$ of combinations include that

$$\binom{n}{n-k} = \binom{n}{k}, \quad (4.3)$$

$$\sum_{k=0}^n \binom{n}{k} = 2^n, \quad (4.4)$$

$$\binom{n-1}{k-1} + \binom{n-1}{k} = \binom{n}{k}, \quad (4.5)$$

$$\binom{n}{k} = \frac{n-k+1}{k} \binom{n}{k-1} \quad (4.6)$$

$$= \frac{k+1}{n-k} \binom{n}{k+1} \quad (4.7)$$

$$= \frac{n}{k} \binom{n-1}{k-1} \quad (4.8)$$

$$= \frac{n}{n-k} \binom{n-1}{k}. \quad (4.9)$$

Equation (4.3) results from changing the variable $k \leftarrow n-k$ in (4.2). Equation (4.4) comes directly from the observation (made at the head of this section) that 2^n total combinations are possible if any k is allowed. Equation (4.5) is seen when an n th block—let us say that it is a black block—is added to an existing set of $n-1$ blocks; to choose k blocks then, you can either choose k from the original set, or the black block plus $k-1$ from the original set. Equations (4.6) through (4.9) come directly from the definition (4.2); they relate combinatoric coefficients to their neighbors in Pascal's triangle (§ 4.2.2).

Because one can choose neither fewer than zero nor more than n from n blocks,

$$\binom{n}{k} = 0 \quad \text{unless } 0 \leq k \leq n. \quad (4.10)$$

For $\binom{n}{k}$ when $n < 0$, there is no obvious definition.

Figure 4.1: The plan for Pascal's triangle.

$$\begin{array}{c}
\binom{0}{0} \\
\binom{1}{0} \binom{1}{1} \\
\binom{2}{0} \binom{2}{1} \binom{2}{2} \\
\binom{3}{0} \binom{3}{1} \binom{3}{2} \binom{3}{3} \\
\binom{4}{0} \binom{4}{1} \binom{4}{2} \binom{4}{3} \binom{4}{4} \\
\vdots
\end{array}$$

4.2.2 Pascal's triangle

Consider the triangular layout in Fig. 4.1 of the various possible $\binom{n}{k}$. Evaluated, this yields Fig. 4.2, *Pascal's triangle*. Notice how each entry in the triangle is the sum of the two entries immediately above, as (4.5) predicts. (In fact this is the easy way to fill Pascal's triangle out: for each entry, just add the two entries above.)

4.3 The binomial theorem

This section presents the binomial theorem and one of its significant consequences.

4.3.1 Expanding the binomial

The *binomial theorem* holds that⁴

$$(a + b)^n = \sum_{k=0}^n \binom{n}{k} a^{n-k} b^k. \quad (4.11)$$

⁴The author is given to understand that, by an heroic derivational effort, (4.11) can be extended to nonintegral n . However, since applied mathematics does not usually concern itself with hard theorems of little known practical use, the extension is not covered in this book.

Figure 4.2: Pascal's triangle.

$$\begin{array}{ccccccc}
& & & & 1 & & \\
& & & 1 & 1 & & \\
& & 1 & 2 & 1 & & \\
& 1 & 3 & 3 & 1 & & \\
1 & 4 & 6 & 4 & 1 & & \\
& 1 & 5 & A & A & 5 & 1 \\
& 1 & 6 & F & 14 & F & 6 & 1 \\
1 & 7 & 15 & 23 & 23 & 15 & 7 & 1 \\
& & & & \vdots & & &
\end{array}$$

In the common case that $a = 1$, $b = \epsilon$, $|\epsilon| \ll 1$, this is

$$(1 + \epsilon)^n = \sum_{k=0}^n \binom{n}{k} \epsilon^k \quad (4.12)$$

(actually this holds for any ϵ , small or large; but the typical case of interest has $|\epsilon| \ll 1$). In either form, the binomial theorem is a direct consequence of the combinatorics of § 4.2. Since

$$(a + b)^n = (a + b)(a + b) \cdots (a + b)(a + b),$$

each $(a+b)$ factor corresponds to one of the “wooden blocks,” where a means rejecting the block and b , accepting it.

4.3.2 Powers of numbers near unity

Since $\binom{n}{0} = 1$ and $\binom{n}{1} = n$, it follows from (4.12) that

$$(1 + \epsilon_o)^n \approx 1 + n\epsilon_o,$$

to excellent precision for nonnegative integral n if ϵ_o is sufficiently small. Furthermore, raising the equation to the $1/n$ power then changing $\delta \leftarrow n\epsilon_o$ and $m \leftarrow n$, we have that

$$(1 + \delta)^{1/m} \approx 1 + \frac{\delta}{m}.$$

Changing $1 + \delta \leftarrow (1 + \epsilon)^n$ and observing from the $(1 + \epsilon_o)^n$ equation above that this implies that $\delta \approx n\epsilon$, we have that

$$(1 + \epsilon)^{n/m} \approx 1 + \frac{n}{m}\epsilon.$$

Inverting this equation yields

$$(1 + \epsilon)^{-n/m} \approx \frac{1}{1 + (n/m)\epsilon} = \frac{[1 - (n/m)\epsilon]}{[1 - (n/m)\epsilon][1 + (n/m)\epsilon]} \approx 1 - \frac{n}{m}\epsilon.$$

Taken together, the last two equations imply that

$$(1 + \epsilon)^x \approx 1 + x\epsilon \tag{4.13}$$

for any real x .

The writer knows of no conventional name⁵ for (4.13), but named or unnamed it is an important equation. The equation offers a simple, accurate way of approximating any real power of numbers in the near neighborhood of 1.

4.3.3 Complex powers of numbers near unity

Equation (4.13) raises the question: what if ϵ or x , or both, are complex? Changing the symbol $z \leftarrow x$ and observing that the infinitesimal ϵ may also be complex, one wants to know whether

$$(1 + \epsilon)^z \approx 1 + z\epsilon \tag{4.14}$$

still holds. No work we have yet done in the book answers the question, because although a complex infinitesimal ϵ poses no particular problem, the action of a complex power z remains undefined. Still, for consistency's sake, one would like (4.14) to hold. In fact nothing prevents us from defining the action of a complex power such that (4.14) does hold, which we now do, logically extending the known result (4.13) into the new domain.

Section 5.4 will investigate the extremely interesting effects which arise when $\Re(\epsilon) = 0$ and the power z in (4.14) grows large, but for the moment we shall use the equation in a more ordinary manner to develop the concept and basic application of the derivative, as follows.

⁵Other than “the first-order Taylor expansion,” but such an unwieldy name does not fit the present context. The Taylor expansion as such will be introduced in Ch. 8.

4.4 The derivative

Having laid down (4.14), we now stand in a position properly to introduce the chapter's subject, the derivative. What is the derivative? The *derivative* is the instantaneous rate or slope of a function. In mathematical symbols and for the moment using real numbers,

$$f'(t) \equiv \lim_{\epsilon \rightarrow 0^+} \frac{f(t + \epsilon/2) - f(t - \epsilon/2)}{\epsilon}. \quad (4.15)$$

Alternately, one can define the same derivative in the unbalanced form

$$f'(t) = \lim_{\epsilon \rightarrow 0^+} \frac{f(t + \epsilon) - f(t)}{\epsilon},$$

but this book generally prefers the more elegant balanced form (4.15), which we will now use in developing the derivative's several properties through the rest of the chapter.⁶

4.4.1 The derivative of the power series

In the very common case that $f(t)$ is the power series

$$f(t) = \sum_{k=-\infty}^{\infty} c_k t^k, \quad (4.16)$$

where the c_k are in general complex coefficients, (4.15) says that

$$\begin{aligned} f'(t) &= \sum_{k=-\infty}^{\infty} \lim_{\epsilon \rightarrow 0^+} \frac{(c_k)(t + \epsilon/2)^k - (c_k)(t - \epsilon/2)^k}{\epsilon} \\ &= \sum_{k=-\infty}^{\infty} \lim_{\epsilon \rightarrow 0^+} c_k t^k \frac{(1 + \epsilon/2t)^k - (1 - \epsilon/2t)^k}{\epsilon}. \end{aligned}$$

Applying (4.14), this is

$$f'(t) = \sum_{k=-\infty}^{\infty} \lim_{\epsilon \rightarrow 0^+} c_k t^k \frac{(1 + k\epsilon/2t) - (1 - k\epsilon/2t)}{\epsilon},$$

⁶From this section through § 4.7, the mathematical notation grows a little thick. There is no helping this. The reader is advised to tread through these sections line by stubborn line, in the good trust that the math thus gained will prove both interesting and useful.

which simplifies to

$$f'(t) = \sum_{k=-\infty}^{\infty} c_k k t^{k-1}. \quad (4.17)$$

Equation (4.17) gives the general derivative of the power series.⁷

4.4.2 The Leibnitz notation

The $f'(t)$ notation used above for the derivative is due to Sir Isaac Newton, and is easier to start with. Usually better on the whole, however, is G.W. Leibnitz's notation⁸

$$\begin{aligned} dt &= \epsilon, \\ df &= f(t + dt/2) - f(t - dt/2), \end{aligned}$$

such that per (4.15),

$$f'(t) = \frac{df}{dt}. \quad (4.18)$$

Here dt is the infinitesimal, and df is a dependent infinitesimal whose size *relative to* dt depends on the independent variable t . For the independent infinitesimal dt , conceptually, one can choose any infinitesimal size ϵ . Usually the exact choice of size does not matter, but occasionally when there are two independent variables it helps the analysis to adjust the size of one of the independent infinitesimals with respect to the other.

The meaning of the symbol d unfortunately depends on the context. In (4.18), the meaning is clear enough: $d(\cdot)$ signifies how much (\cdot) changes when the independent variable t increments by dt .⁹ Notice, however, that the notation dt itself has two distinct meanings:¹⁰

⁷Equation (4.17) admittedly has not explicitly considered what happens when the real t becomes the complex z , but § 4.4.3 will remedy the oversight.

⁸This subsection is likely to confuse many readers the first time they read it. The reason is that Leibnitz elements like dt and ∂f usually tend to appear in practice in certain specific relations to one another, like $\partial f / \partial z$. As a result, many users of applied mathematics have never developed a clear understanding as to precisely what the individual symbols mean. Often they have developed positive misunderstandings. Because there is significant practical benefit in learning how to handle the Leibnitz notation correctly—particularly in applied complex variable theory—this subsection seeks to present each Leibnitz element in its correct light.

⁹If you do not fully understand this sentence, reread it carefully with reference to (4.15) and (4.18) until you do; it's important.

¹⁰This is difficult, yet the author can think of no clearer, more concise way to state it. The quantities dt and df represent coordinated infinitesimal changes in t and f respectively,

- the independent infinitesimal $dt = \epsilon$; and
- $d(t)$, which is how much (t) changes as t increments by dt .

At first glance, the distinction between dt and $d(t)$ seems a distinction without a difference; and for most practical cases of interest, so indeed it is. However, when switching perspective in mid-analysis as to which variables are dependent and which are independent, or when changing multiple independent complex variables simultaneously, the math can get a little tricky.

so there is usually no trouble with treating dt and df as though they were the same kind of thing. However, at the fundamental level they really aren't.

If t is an independent variable, then dt is just an infinitesimal of some kind, whose specific size could be a function of t but more likely is just a constant. If a constant, then dt does not fundamentally have anything to do with t as such. In fact, if s and t are both independent variables, then we can (and in complex analysis sometimes do) say that $ds = dt = \epsilon$, after which nothing prevents us from using the symbols ds and dt interchangeably. Maybe it would be clearer in some cases to write ϵ instead of dt , but the latter is how it is conventionally written.

By contrast, if f is a dependent variable, then df or $d(f)$ is the amount by which f changes as t changes by dt . The df is infinitesimal but not constant; it is a function of t . Maybe it would be clearer in some cases to write $d_t f$ instead of df , but for most cases the former notation is unnecessarily cluttered; the latter is how it is conventionally written.

Now, most of the time, what we are interested in is not dt or df as such, but rather the ratio df/dt or the sum $\sum_k f(k dt) dt = \int f(t) dt$. For this reason, we do not usually worry about which of df and dt is the independent infinitesimal, nor do we usually worry about the precise value of dt . This leads one to forget that dt does indeed have a precise value. What confuses is when one changes perspective in mid-analysis, now regarding f as the independent variable. Changing perspective is allowed and perfectly proper, but one must take care: the dt and df after the change are not the same as the dt and df before the change. However, the ratio df/dt remains the same in any case.

Sometimes when writing a differential equation like the potential-kinetic energy equation $ma dx = mv dv$, we do not necessarily have either v or x in mind as the independent variable. This is fine. The important point is that dv and dx be coordinated so that the ratio dv/dx has a definite value no matter which of the two be regarded as independent, or whether the independent be some third variable (like t) not in the equation.

One can avoid the confusion simply by keeping the dv/dx or df/dt always in ratio, never treating the infinitesimals individually. Many applied mathematicians do precisely that. That is okay as far as it goes, but it really denies the entire point of the Leibnitz notation. One might as well just stay with the Newton notation in that case. Instead, this writer recommends that you learn the Leibnitz notation properly, developing the ability to treat the infinitesimals individually.

Because the book is a book of applied mathematics, this footnote does not attempt to say everything there is to say about infinitesimals. For instance, it has not yet pointed out (but does so now) that even if s and t are equally independent variables, one can have $dt = \epsilon(t)$, $ds = \delta(s, t)$, such that dt has prior independence to ds . The point is not to fathom all the possible implications from the start; you can do that as the need arises. The point is to develop a clear picture in your mind of what a Leibnitz infinitesimal really is. Once you have the picture, you can go from there.

In such cases, it may be wise to use the symbol dt to mean $d(t)$ only, introducing some unambiguous symbol like ϵ to represent the independent infinitesimal. In any case you should appreciate the conceptual difference between $dt = \epsilon$ and $d(t)$.

Where two or more independent variables are at work in the same equation, it is conventional to use the symbol ∂ instead of d , as a reminder that the reader needs to pay attention to which ∂ tracks which independent variable.¹¹ (If needed or desired, one can write $\partial_t(\cdot)$ when tracking t , $\partial_s(\cdot)$ when tracking s , etc. Use discretion, though. Such notation appears only rarely in the literature, so your audience might not understand it when you write it.) Conventional shorthand for $d(df)$ is d^2f ; for $(dt)^2$, dt^2 ; so

$$\frac{d(df/dt)}{dt} = \frac{d^2f}{dt^2}$$

is a derivative of a derivative, or *second derivative*. By extension, the notation

$$\frac{d^k f}{dt^k}$$

represents the k th derivative.

4.4.3 The derivative of a function of a complex variable

For (4.15) to be robust, written here in the slightly more general form

$$\frac{df}{dz} = \lim_{\epsilon \rightarrow 0} \frac{f(z + \epsilon/2) - f(z - \epsilon/2)}{\epsilon}, \quad (4.19)$$

one should like it to evaluate the same in the limit regardless of the complex phase of ϵ . That is, if δ is a positive real infinitesimal, then it should be equally valid to let $\epsilon = \delta$, $\epsilon = -\delta$, $\epsilon = i\delta$, $\epsilon = -i\delta$, $\epsilon = (4 - i3)\delta$ or any other infinitesimal value, so long as $0 < |\epsilon| \ll 1$. One should like the derivative (4.19) to come out the same regardless of the Argand direction from which ϵ approaches 0 (see Fig. 2.5). In fact for the sake of robustness, one normally demands that derivatives do come out the same regardless of the Argand direction; and (4.19) rather than (4.15) is the definition we normally use for the derivative for this reason. Where the limit (4.19) is sensitive to the Argand direction or complex phase of ϵ , there we normally say that the derivative does not exist.

¹¹The writer confesses that he remains unsure why this minor distinction merits the separate symbol ∂ , but he accepts the notation as conventional nevertheless.

Where the derivative (4.19) does exist—where the derivative is finite and insensitive to Argand direction—there we say that the function $f(z)$ is *differentiable*.

Excepting the nonanalytic parts of complex numbers ($|\cdot|$, $\arg[\cdot]$, $[\cdot]^*$, $\Re[\cdot]$ and $\Im[\cdot]$; see § 2.12.3), plus the Heaviside unit step $u(t)$ and the Dirac delta $\delta(t)$ (§ 7.7), most functions encountered in applications do meet the criterion (4.19) except at isolated nonanalytic points (like $z = 0$ in $h[z] = 1/z$ or $g[z] = \sqrt{z}$). Meeting the criterion, such functions are fully differentiable except at their poles (where the derivative goes infinite in any case) and other nonanalytic points. Particularly, the key formula (4.14), written here as

$$(1 + \epsilon)^w \approx 1 + w\epsilon,$$

works without modification when ϵ is complex; so the derivative (4.17) of the general power series,

$$\frac{d}{dz} \sum_{k=-\infty}^{\infty} c_k z^k = \sum_{k=-\infty}^{\infty} c_k k z^{k-1} \quad (4.20)$$

holds equally well for complex z as for real.

4.4.4 The derivative of z^a

Inspection of § 4.4.1's logic in light of (4.14) reveals that nothing prevents us from replacing the real t , real ϵ and integral k of that section with arbitrary complex z , ϵ and a . That is,

$$\begin{aligned} \frac{d(z^a)}{dz} &= \lim_{\epsilon \rightarrow 0} \frac{(z + \epsilon/2)^a - (z - \epsilon/2)^a}{\epsilon} \\ &= \lim_{\epsilon \rightarrow 0} z^a \frac{(1 + \epsilon/2z)^a - (1 - \epsilon/2z)^a}{\epsilon} \\ &= \lim_{\epsilon \rightarrow 0} z^a \frac{(1 + a\epsilon/2z) - (1 - a\epsilon/2z)}{\epsilon}, \end{aligned}$$

which simplifies to

$$\frac{d(z^a)}{dz} = az^{a-1} \quad (4.21)$$

for any complex z and a .

How exactly to evaluate z^a or z^{a-1} when a is complex is another matter, treated in § 5.4 and its (5.12); but in any case you can use (4.21) for real a right now.

4.4.5 The logarithmic derivative

Sometimes one is more interested in knowing the rate of $f(t)$ *relative to the value of $f(t)$* than in knowing the absolute rate itself. For example, if you inform me that you earn \$ 1000 a year on a bond you hold, then I may commend you vaguely for your thrift but otherwise the information does not tell me much. However, if you inform me instead that you earn 10 percent a year on the same bond, then I might want to invest. The latter figure is a relative rate, or *logarithmic derivative*,

$$\frac{df/dt}{f(t)} = \frac{d}{dt} \ln f(t). \quad (4.22)$$

The investment principal grows at the absolute rate df/dt , but the bond's interest rate is $(df/dt)/f(t)$.

The natural logarithmic notation $\ln f(t)$ may not mean much to you yet, as we'll not introduce it formally until § 5.2, so you can ignore the right side of (4.22) for the moment; but the equation's left side at least should make sense to you. It expresses the significant concept of a relative rate, like 10 percent annual interest on a bond.

4.5 Basic manipulation of the derivative

This section introduces the derivative chain and product rules.

4.5.1 The derivative chain rule

If f is a function of w , which itself is a function of z , then¹²

$$\frac{df}{dz} = \left(\frac{df}{dw} \right) \left(\frac{dw}{dz} \right). \quad (4.23)$$

Equation (4.23) is the *derivative chain rule*.¹³

4.5.2 The derivative product rule

In general per (4.19),

$$d \left[\prod_j f_j(z) \right] = \prod_j f_j \left(z + \frac{dz}{2} \right) - \prod_j f_j \left(z - \frac{dz}{2} \right).$$

But to first order,

$$f_j \left(z \pm \frac{dz}{2} \right) \approx f_j(z) \pm \left(\frac{df_j}{dz} \right) \left(\frac{dz}{2} \right) = f_j(z) \pm \frac{df_j}{2},$$

so in the limit,

$$d \left[\prod_j f_j(z) \right] = \prod_j \left(f_j(z) + \frac{df_j}{2} \right) - \prod_j \left(f_j(z) - \frac{df_j}{2} \right).$$

¹²For example, one can rewrite

$$f(z) = \sqrt{3z^2 - 1}$$

in the form

$$\begin{aligned} f(w) &= w^{1/2}, \\ w(z) &= 3z^2 - 1. \end{aligned}$$

Then

$$\begin{aligned} \frac{df}{dw} &= \frac{1}{2w^{1/2}} = \frac{1}{2\sqrt{3z^2 - 1}}, \\ \frac{dw}{dz} &= 6z, \end{aligned}$$

so by (4.23),

$$\frac{df}{dz} = \left(\frac{df}{dw} \right) \left(\frac{dw}{dz} \right) = \frac{6z}{2\sqrt{3z^2 - 1}} = \frac{3z}{\sqrt{3z^2 - 1}}.$$

¹³It bears emphasizing to readers who may inadvertently have picked up unhelpful ideas about the Leibnitz notation in the past: the dw factor in the denominator cancels the dw factor in the numerator; a thing divided by itself is 1. That's it. There is nothing more to the proof of the derivative chain rule than that.

Since the product of two or more df_j is negligible compared to the first-order infinitesimals to which they are added here, this simplifies to

$$d \left[\prod_j f_j(z) \right] = \left[\prod_j f_j(z) \right] \left[\sum_k \frac{df_k}{f_k(z)} \right] - \left[\prod_j f_j(z) \right] \left[\sum_k \frac{-df_k}{f_k(z)} \right],$$

or in other words

$$d \prod_j f_j = \left[\prod_j f_j \right] \left[\sum_k \frac{df_k}{f_k} \right]. \quad (4.24)$$

In the common case of only two f_j , this comes to

$$d(f_1 f_2) = f_2 df_1 + f_1 df_2. \quad (4.25)$$

On the other hand, if $f_1(z) = f(z)$ and $f_2(z) = 1/g(z)$, then by the derivative chain rule (4.23), $df_2 = -dg/g^2$; so,

$$d \left(\frac{f}{g} \right) = \frac{g df - f dg}{g^2}. \quad (4.26)$$

Equation (4.24) is the *derivative product rule*.

After studying the complex exponential in Ch. 5, we shall stand in a position to write (4.24) in the slightly specialized but often useful form¹⁴

$$\begin{aligned} & d \left[\prod_j g_j^{a_j} \prod_j e^{b_j h_j} \prod_j \ln c_j p_j \right] \\ &= \left[\prod_j g_j^{a_j} \prod_j e^{b_j h_j} \prod_j \ln c_j p_j \right] \\ &\quad \times \left[\sum_k a_k \frac{dg_k}{g_k} + \sum_k b_k dh_k + \sum_k \frac{dp_k}{p_k \ln c_k p_k} \right]. \end{aligned} \quad (4.27)$$

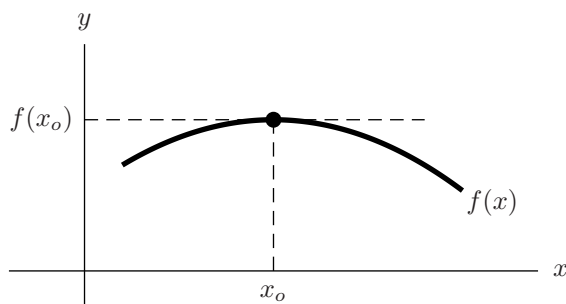
where the a_k , b_k and c_k are arbitrary complex coefficients and the g_k , h_k and p_k are arbitrary functions.¹⁵

¹⁴This paragraph is extra. You can skip it for now if you prefer.

¹⁵The subsection is sufficiently abstract that it is a little hard to understand unless one already knows what it means. An example may help:

$$d \left[\frac{u^2 v^3}{z} e^{-5t} \ln 7s \right] = \left[\frac{u^2 v^3}{z} e^{-5t} \ln 7s \right] \left[2 \frac{du}{u} + 3 \frac{dv}{v} - \frac{dz}{z} - 5 dt + \frac{ds}{s \ln 7s} \right].$$

Figure 4.3: A local extremum.



4.6 Extrema and higher derivatives

One problem which arises very frequently in applied mathematics is the problem of finding a local *extremum*—that is, a local minimum or maximum—of a real-valued function $f(x)$. Refer to Fig. 4.3. The almost distinctive characteristic of the extremum x_o is that¹⁶

$$\left. \frac{df}{dx} \right|_{x=x_o} = 0. \quad (4.28)$$

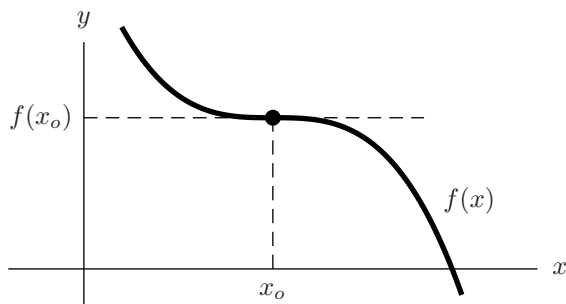
At the extremum, the slope is zero. The curve momentarily runs level there. One solves (4.28) to find the extremum.

Whether the extremum be a minimum or a maximum depends on whether the curve turn from a downward slope to an upward, or from an upward slope to a downward, respectively. If from downward to upward, then the derivative of the slope is evidently positive; if from upward to downward, then negative. But the derivative of the slope is just the derivative of the derivative, or second derivative. Hence if $df/dx = 0$ at $x = x_o$, then

$$\begin{aligned} \left. \frac{d^2 f}{dx^2} \right|_{x=x_o} > 0 & \text{ implies a local minimum at } x_o; \\ \left. \frac{d^2 f}{dx^2} \right|_{x=x_o} < 0 & \text{ implies a local maximum at } x_o. \end{aligned}$$

¹⁶The notation $P|_Q$ means “ P when Q ,” “ P , given Q ,” or “ P evaluated at Q .” Sometimes it is alternately written $P|Q$ or $[P]_Q$.

Figure 4.4: A level inflection.



Regarding the case

$$\left. \frac{d^2 f}{dx^2} \right|_{x=x_o} = 0,$$

this might be either a minimum or a maximum but more probably is neither, being rather a *level inflection point* as depicted in Fig. 4.4.¹⁷ (In general the term *inflection point* signifies a point at which the second derivative is zero. The inflection point of Fig. 4.4 is *level* because its first derivative is zero, too.)

4.7 L'Hôpital's rule

If $z = z_o$ is a root of both $f(z)$ and $g(z)$, or alternately if $z = z_o$ is a pole of both functions—that is, if both functions go to zero or infinity together at $z = z_o$ —then *l'Hôpital's rule* holds that

$$\lim_{z \rightarrow z_o} \frac{f(z)}{g(z)} = \left. \frac{df/dz}{dg/dz} \right|_{z=z_o}. \quad (4.29)$$

¹⁷Of course if the first and second derivatives are zero not just at $x = x_o$ but everywhere, then $f(x) = y_o$ is just a level straight line, but you knew that already. Whether one chooses to call some random point on a level straight line an inflection point or an extremum, or both or neither, would be a matter of definition, best established not by prescription but rather by the needs of the model at hand.

In the case where $z = z_o$ is a root, l'Hôpital's rule is proved by reasoning¹⁸

$$\begin{aligned}\lim_{z \rightarrow z_o} \frac{f(z)}{g(z)} &= \lim_{z \rightarrow z_o} \frac{f(z) - 0}{g(z) - 0} \\ &= \lim_{z \rightarrow z_o} \frac{f(z) - f(z_o)}{g(z) - g(z_o)} = \lim_{z \rightarrow z_o} \frac{df}{dg} = \lim_{z \rightarrow z_o} \frac{df/dz}{dg/dz}.\end{aligned}$$

In the case where $z = z_o$ is a pole, new functions $F(z) \equiv 1/f(z)$ and $G(z) \equiv 1/g(z)$ of which $z = z_o$ is a root are defined, with which

$$\lim_{z \rightarrow z_o} \frac{f(z)}{g(z)} = \lim_{z \rightarrow z_o} \frac{G(z)}{F(z)} = \lim_{z \rightarrow z_o} \frac{dG}{dF} = \lim_{z \rightarrow z_o} \frac{-dg/g^2}{-df/f^2},$$

where we have used the fact from (4.21) that $d(1/u) = -du/u^2$ for any u . Canceling the minus signs and multiplying by g^2/f^2 , we have that

$$\lim_{z \rightarrow z_o} \frac{g(z)}{f(z)} = \lim_{z \rightarrow z_o} \frac{dg}{df}.$$

Inverting,

$$\lim_{z \rightarrow z_o} \frac{f(z)}{g(z)} = \lim_{z \rightarrow z_o} \frac{df}{dg} = \lim_{z \rightarrow z_o} \frac{df/dz}{dg/dz}.$$

And if z_o itself is infinite? Then, whether it represents a root or a pole, we define the new variable $Z = 1/z$ and the new functions $\Phi(Z) = f(1/Z) = f(z)$ and $\Gamma(Z) = g(1/Z) = g(z)$, with which we apply l'Hôpital's rule for $Z \rightarrow 0$ to obtain

$$\begin{aligned}\lim_{z \rightarrow \infty} \frac{f(z)}{g(z)} &= \lim_{Z \rightarrow 0} \frac{\Phi(Z)}{\Gamma(Z)} = \lim_{Z \rightarrow 0} \frac{d\Phi/dZ}{d\Gamma/dZ} = \lim_{Z \rightarrow 0} \frac{df/dZ}{dg/dZ} \\ &= \lim_{\substack{z \rightarrow \infty \\ Z \rightarrow 0}} \frac{(df/dz)(dz/dZ)}{(dg/dz)(dz/dZ)} = \lim_{z \rightarrow \infty} \frac{(df/dz)(-z^2)}{(dg/dz)(-z^2)} = \lim_{z \rightarrow \infty} \frac{df/dz}{dg/dz}.\end{aligned}$$

Nothing in the derivation requires that z or z_o be real. Nothing prevents one from applying l'Hôpital's rule recursively, should the occasion arise.¹⁹

¹⁸Partly with reference to [46, "L'Hopital's rule," 03:40, 5 April 2006].

¹⁹Consider for example the ratio $\lim_{x \rightarrow 0} (x^3 + x)^2/x^2$, which is $0/0$. The easier way to resolve this particular ratio would naturally be to cancel a factor of x^2 from it; but just to make the point let us apply l'Hôpital's rule instead, reducing the ratio to $\lim_{x \rightarrow 0} 2(x^3 + x)(3x^2 + 1)/2x$, which is still $0/0$. Applying l'Hôpital's rule again to the result yields $\lim_{x \rightarrow 0} 2[(3x^2 + 1)^2 + (x^3 + x)(6x)]/2 = 2/2 = 1$. Where expressions involving trigonometric or special functions (Chs. 3, 5 and [not yet written]) appear in ratio, a recursive application of l'Hôpital's rule can be just the thing one needs.

Observe that one must stop applying l'Hôpital's rule once the ratio is no longer $0/0$ or ∞/∞ . In the example, applying the rule a third time would have ruined the result.

L'Hôpital's rule is used in evaluating indeterminate forms of the kinds $0/0$ and ∞/∞ , plus related forms like $(0)(\infty)$ which can be recast into either of the two main forms. Good examples of the use require math from Ch. 5 and later, but if we may borrow from (5.7) the natural logarithmic function and its derivative,²⁰

$$\frac{d}{dx} \ln x = \frac{1}{x},$$

then a typical l'Hôpital example is²¹

$$\lim_{x \rightarrow \infty} \frac{\ln x}{\sqrt{x}} = \lim_{x \rightarrow \infty} \frac{1/x}{1/2\sqrt{x}} = \lim_{x \rightarrow \infty} \frac{2}{\sqrt{x}} = 0.$$

The example incidentally shows that natural logarithms grow slower than square roots, an instance of a more general principle we shall meet in § 5.3.

Section 5.3 will put l'Hôpital's rule to work.

4.8 The Newton-Raphson iteration

The *Newton-Raphson iteration* is a powerful, fast converging, broadly applicable method for finding roots numerically. Given a function $f(z)$ of which the root is desired, the Newton-Raphson iteration is

$$z_{k+1} = z - \left. \frac{f(z)}{\frac{d}{dz}f(z)} \right|_{z=z_k}. \quad (4.30)$$

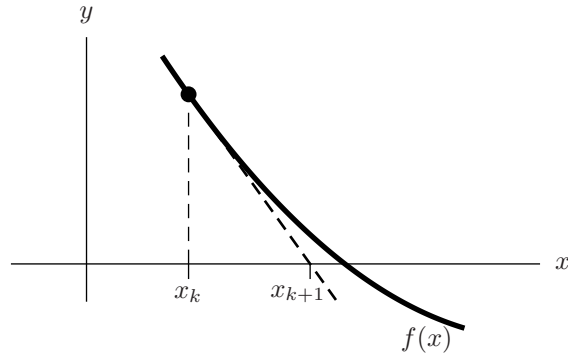
One begins the iteration by guessing the root and calling the guess z_0 . Then z_1 , z_2 , z_3 , etc., calculated in turn by the iteration (4.30), give successively better estimates of the true root z_∞ .

To understand the Newton-Raphson iteration, consider the function $y = f(x)$ of Fig 4.5. The iteration approximates the curve $f(x)$ by its tangent

²⁰This paragraph is optional reading for the moment. You can read Ch. 5 first, then come back here and read the paragraph if you prefer.

²¹[37, § 10-2]

Figure 4.5: The Newton-Raphson iteration.



line²² (shown as the dashed line in the figure):

$$\tilde{f}_k(x) = f(x_k) + \left[\frac{d}{dx} f(x) \right]_{x=x_k} (x - x_k).$$

It then approximates the root x_{k+1} as the point at which $\tilde{f}_k(x_{k+1}) = 0$:

$$\tilde{f}_k(x_{k+1}) = 0 = f(x_k) + \left[\frac{d}{dx} f(x) \right]_{x=x_k} (x_{k+1} - x_k).$$

Solving for x_{k+1} , we have that

$$x_{k+1} = x - \frac{f(x)}{\frac{d}{dx} f(x)} \Big|_{x=x_k},$$

which is (4.30) with $x \leftarrow z$.

Although the illustration uses real numbers, nothing forbids complex z and $f(z)$. The Newton-Raphson iteration works just as well for these.

²²A *tangent* line, also just called a *tangent*, is the line which most nearly approximates a curve at a given point. The tangent touches the curve at the point, and in the neighborhood of the point it goes in the same direction the curve goes. The dashed line in Fig. 4.5 is a good example of a tangent line.

The relationship between the tangent line and the trigonometric tangent function of Ch. 3 is slightly obscure, maybe more of linguistic interest than of mathematical. The trigonometric tangent function is named from a variation on Fig. 3.1 in which the triangle's bottom leg is extended to unit length, leaving the rightward leg tangent to the circle.

The principal limitation of the Newton-Raphson arises when the function has more than one root, as most interesting functions do. The iteration often converges on the root nearest the initial guess z_0 but does not always, and in any case there is no guarantee that the root it finds is the one you wanted. The most straightforward way to beat this problem is to find *all* the roots: first you find some root α , then you remove that root (without affecting any of the other roots) by dividing $f(z)/(z - \alpha)$, then you find the next root by iterating on the new function $f(z)/(z - \alpha)$, and so on until you have found all the roots. If this procedure is not practical (perhaps because the function has a large or infinite number of roots), then you should probably take care to make a sufficiently accurate initial guess if you can.

A second limitation of the Newton-Raphson is that, if you happen to guess z_0 especially unfortunately, then the iteration might never converge at all. For example, the roots of $f(z) = z^2 + 2$ are $z = \pm i\sqrt{2}$, but if you guess that $z_0 = 1$ then the iteration has no way to leave the real number line, so it never converges²³ (and if you guess that $z_0 = \sqrt{2}$ —well, try it with your pencil and see what z_2 comes out to be). You can fix the problem with a different, possibly complex initial guess.

A third limitation arises where there is a multiple root. In this case, the Newton-Raphson normally still converges, but relatively slowly. For instance, the Newton-Raphson converges relatively slowly on the triple root of $f(z) = z^3$. However, even the relatively slow convergence is still pretty fast and is usually adequate, even for calculations by hand.

Usually in practice, the Newton-Raphson iteration works very well. For most functions, once the Newton-Raphson finds the root's neighborhood, it converges on the actual root remarkably quickly. Figure 4.5 shows why: in the neighborhood, the curve hardly departs from the straight line.

The Newton-Raphson iteration is a champion square root calculator, incidentally. Consider

$$f(x) = x^2 - p,$$

whose roots are

$$x = \pm\sqrt{p}.$$

Per (4.30), the Newton-Raphson iteration for this is

$$x_{k+1} = \frac{1}{2} \left[x_k + \frac{p}{x_k} \right]. \quad (4.31)$$

If you start by guessing

$$x_0 = 1$$

²³It is entertaining to try this on a computer. Then try again with $z_0 = 1 + i2^{-0 \times 10}$.

and iterate several times, the iteration (4.31) converges on $x_\infty = \sqrt[n]{p}$ fast. To calculate the n th root $x = p^{1/n}$, let

$$f(x) = x^n - p$$

and iterate^{24,25}

$$x_{k+1} = \frac{1}{n} \left[(n-1)x_k + \frac{p}{x_k^{n-1}} \right]. \quad (4.32)$$

Section 13.7 generalizes the Newton-Raphson iteration to handle vector-valued functions.

This concludes the chapter. Chapter 8, treating the Taylor series, will continue the general discussion of the derivative.

²⁴Equations (4.31) and (4.32) work not only for real p but also usually for complex. Given $x_0 = 1$, however, they converge reliably and orderly only for real, nonnegative p . (To see why, sketch $f[x]$ in the fashion of Fig. 4.5.)

If reliable, orderly convergence is needed for complex $p = u + iv = \sigma \operatorname{cis} \psi$, $\sigma \geq 0$, you can decompose $p^{1/n}$ per de Moivre's theorem (3.26) as $p^{1/n} = \sigma^{1/n} \operatorname{cis}(\psi/n)$, in which $\operatorname{cis}(\psi/n) = \cos(\psi/n) + i \sin(\psi/n)$ is calculated by the Taylor series of Table 8.1. Then σ is real and nonnegative, upon which (4.32) reliably, orderly computes $\sigma^{1/n}$.

The Newton-Raphson iteration however excels as a *practical* root-finding technique, so it often pays to be a little less theoretically rigid in applying it. If so, then don't bother to decompose; seek $p^{1/n}$ directly, using complex z_k in place of the real x_k . In the uncommon event that the direct iteration does not seem to converge, start over again with some randomly chosen complex z_0 . This saves effort and usually works.

²⁵[37, § 4-9][30, § 6.1.1][45]

Chapter 5

The complex exponential

In higher mathematics, the complex natural exponential is almost impossible to avoid. It seems to appear just about everywhere. This chapter introduces the concept of the natural exponential and of its inverse, the natural logarithm; and shows how the two operate on complex arguments. It derives the functions' basic properties and explains their close relationship to the trigonometrics. It works out the functions' derivatives and the derivatives of the basic members of the trigonometric and inverse trigonometric families to which they respectively belong.

5.1 The real exponential

Consider the factor

$$(1 + \epsilon)^N.$$

This is the overall factor by which a quantity grows after N iterative rounds of multiplication by $(1 + \epsilon)$. What happens when ϵ is very small but N is very large? The really interesting question is, what happens in the limit, as $\epsilon \rightarrow 0$ and $N \rightarrow \infty$, while $x = \epsilon N$ remains a finite number? The answer is that the factor becomes

$$\exp x \equiv \lim_{\epsilon \rightarrow 0} (1 + \epsilon)^{x/\epsilon}. \quad (5.1)$$

Equation (5.1) defines the *natural exponential function*—commonly, more briefly named the *exponential function*. Another way to write the same definition is

$$\exp x = e^x, \quad (5.2)$$

$$e \equiv \lim_{\epsilon \rightarrow 0} (1 + \epsilon)^{1/\epsilon}. \quad (5.3)$$

Whichever form we write it in, the question remains as to whether the limit actually exists; that is, whether $0 < e < \infty$; whether in fact we can put some concrete bound on e . To show that we can,¹ we observe per (4.19) that the derivative of the exponential function is

$$\begin{aligned}
 \frac{d}{dx} \exp x &= \lim_{\delta \rightarrow 0} \frac{\exp(x + \delta/2) - \exp(x - \delta/2)}{\delta} \\
 &= \lim_{\delta, \epsilon \rightarrow 0} \frac{(1 + \epsilon)^{(x+\delta/2)/\epsilon} - (1 + \epsilon)^{(x-\delta/2)/\epsilon}}{\delta} \\
 &= \lim_{\delta, \epsilon \rightarrow 0} (1 + \epsilon)^{x/\epsilon} \frac{(1 + \epsilon)^{+\delta/2\epsilon} - (1 + \epsilon)^{-\delta/2\epsilon}}{\delta} \\
 &= \lim_{\delta, \epsilon \rightarrow 0} (1 + \epsilon)^{x/\epsilon} \frac{(1 + \delta/2) - (1 - \delta/2)}{\delta} \\
 &= \lim_{\epsilon \rightarrow 0} (1 + \epsilon)^{x/\epsilon},
 \end{aligned}$$

which is to say that

$$\frac{d}{dx} \exp x = \exp x. \quad (5.4)$$

This is a curious, important result: the derivative of the exponential function is the exponential function itself; the slope and height of the exponential function are everywhere equal. For the moment, however, what interests us is that

$$\frac{d}{dx} \exp 0 = \exp 0 = \lim_{\epsilon \rightarrow 0} (1 + \epsilon)^0 = 1,$$

which says that the slope and height of the exponential function are both unity at $x = 0$, implying that the straight line which best approximates the exponential function in that neighborhood—the *tangent line*, which just grazes the curve—is

$$y(x) = 1 + x.$$

With the tangent line $y(x)$ found, the next step toward putting a concrete bound on e is to show that $y(x) \leq \exp x$ for all real x , that the curve runs nowhere below the line. To show this, we observe per (5.1) that the essential action of the exponential function is to multiply repeatedly by $1 + \epsilon$ as x increases, to divide repeatedly by $1 + \epsilon$ as x decreases. Since $1 + \epsilon > 1$, this

¹Excepting (5.4), the author would prefer to omit much of the rest of this section, but even at the applied level cannot think of a logically permissible way to do it. It seems nonobvious that the limit $\lim_{\epsilon \rightarrow 0} (1 + \epsilon)^{1/\epsilon}$ actually does exist. The rest of this section shows why it does.

action means for real x that

$$\exp x_1 \leq \exp x_2 \quad \text{if } x_1 \leq x_2.$$

However, a positive number remains positive no matter how many times one multiplies or divides it by $1 + \epsilon$, so the same action also means that

$$0 \leq \exp x$$

for all real x . In light of (5.4), the last two equations imply further that

$$\begin{aligned} \frac{d}{dx} \exp x_1 &\leq \frac{d}{dx} \exp x_2 \quad \text{if } x_1 \leq x_2, \\ 0 &\leq \frac{d}{dx} \exp x. \end{aligned}$$

But we have purposely defined the tangent line $y(x) = 1 + x$ such that

$$\begin{aligned} \exp 0 &= y(0) = 1, \\ \frac{d}{dx} \exp 0 &= \frac{d}{dx} y(0) = 1; \end{aligned}$$

that is, such that the line just grazes the curve of $\exp x$ at $x = 0$. Rightward, at $x > 0$, evidently the curve's slope only increases, bending upward away from the line. Leftward, at $x < 0$, evidently the curve's slope only decreases, again bending upward away from the line. Either way, the curve never crosses below the line for real x . In symbols,

$$y(x) \leq \exp x.$$

Figure 5.1 depicts.

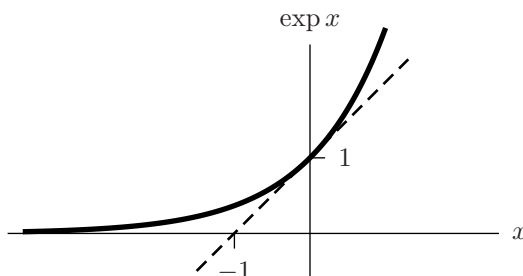
Evaluating the last inequality at $x = -1/2$ and $x = 1$, we have that

$$\begin{aligned} \frac{1}{2} &\leq \exp\left(-\frac{1}{2}\right), \\ 2 &\leq \exp(1). \end{aligned}$$

But per (5.2), $\exp x = e^x$, so

$$\begin{aligned} \frac{1}{2} &\leq e^{-1/2}, \\ 2 &\leq e^1, \end{aligned}$$

Figure 5.1: The natural exponential.



or in other words,

$$2 \leq e \leq 4, \quad (5.5)$$

which in consideration of (5.2) puts the desired bound on the exponential function. The limit does exist.

By the Taylor series of Table 8.1, the value

$$e \approx 0x2.B7E1$$

can readily be calculated, but the derivation of that series does not come until Ch. 8.

5.2 The natural logarithm

In the general exponential expression b^x , one can choose any base b ; for example, $b = 2$ is an interesting choice. As we shall see in § 5.4, however, it turns out that $b = e$, where e is the constant introduced in (5.3), is the most interesting choice of all. For this reason among others, the base- e logarithm is similarly interesting, such that we define for it the special notation

$$\ln(\cdot) = \log_e(\cdot),$$

and call it the *natural logarithm*. Just as for any other base b , so also for base $b = e$; thus the natural logarithm inverts the natural exponential and vice versa:

$$\begin{aligned} \ln \exp x &= \ln e^x = x, \\ \exp \ln x &= e^{\ln x} = x. \end{aligned} \quad (5.6)$$

Figure 5.2: The natural logarithm.

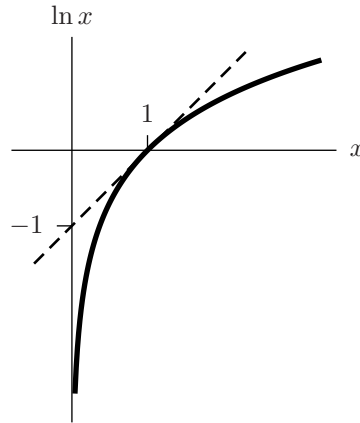


Figure 5.2 plots the natural logarithm.

If

$$y = \ln x,$$

then

$$x = \exp y,$$

and per (5.4),

$$\frac{dx}{dy} = \exp y.$$

But this means that

$$\frac{dx}{dy} = x,$$

the inverse of which is

$$\frac{dy}{dx} = \frac{1}{x}.$$

In other words,

$$\frac{d}{dx} \ln x = \frac{1}{x}. \quad (5.7)$$

Like many of the equations in these early chapters, here is another rather significant result.²

²Besides the result itself, the technique which leads to the result is also interesting and is worth mastering. We shall use the technique more than once in this book.

One can specialize Table 2.5's logarithmic base-conversion identity to read

$$\log_b w = \frac{\ln w}{\ln b}. \quad (5.8)$$

This equation converts any logarithm to a natural logarithm. Base $b = 2$ logarithms are interesting, so we note here that

$$\ln 2 = -\ln \frac{1}{2} \approx 0.693147,$$

which Ch. 8 and its Table 8.1 will show how to calculate.

5.3 Fast and slow functions

The exponential $\exp x$ is a *fast function*. The logarithm $\ln x$ is a *slow function*. These functions grow, diverge or decay respectively faster and slower than x^a .

Such claims are proved by l'Hôpital's rule (4.29). Applying the rule, we have that

$$\begin{aligned} \lim_{x \rightarrow \infty} \frac{\ln x}{x^a} &= \lim_{x \rightarrow \infty} \frac{-1}{ax^a} = \begin{cases} 0 & \text{if } a > 0, \\ +\infty & \text{if } a \leq 0, \end{cases} \\ \lim_{x \rightarrow 0} \frac{\ln x}{x^a} &= \lim_{x \rightarrow 0} \frac{-1}{ax^a} = \begin{cases} -\infty & \text{if } a \geq 0, \\ 0 & \text{if } a < 0, \end{cases} \end{aligned} \quad (5.9)$$

which reveals the logarithm to be a slow function. Since the $\exp(\cdot)$ and $\ln(\cdot)$ functions are mutual inverses, we can leverage (5.9) to show also that

$$\begin{aligned} \lim_{x \rightarrow \infty} \frac{\exp(\pm x)}{x^a} &= \lim_{x \rightarrow \infty} \exp \left[\ln \frac{\exp(\pm x)}{x^a} \right] \\ &= \lim_{x \rightarrow \infty} \exp [\pm x - a \ln x] \\ &= \lim_{x \rightarrow \infty} \exp \left[(x) \left(\pm 1 - a \frac{\ln x}{x} \right) \right] \\ &= \lim_{x \rightarrow \infty} \exp [(x) (\pm 1 - 0)] \\ &= \lim_{x \rightarrow \infty} \exp [\pm x]. \end{aligned}$$

That is,

$$\begin{aligned} \lim_{x \rightarrow \infty} \frac{\exp(+x)}{x^a} &= \infty, \\ \lim_{x \rightarrow \infty} \frac{\exp(-x)}{x^a} &= 0, \end{aligned} \quad (5.10)$$

which reveals the exponential to be a fast function. Exponentials grow or decay faster than powers; logarithms diverge slower.

Such conclusions are extended to bases other than the natural base e simply by observing that $\log_b x = \ln x / \ln b$ and that $b^x = \exp(x \ln b)$. Thus exponentials generally are fast and logarithms generally are slow, regardless of the base.³

It is interesting and worthwhile to contrast the sequence

$$\cdots, -\frac{3!}{x^4}, \frac{2!}{x^3}, -\frac{1!}{x^2}, \frac{0!}{x^1}, \frac{x^0}{0!}, \frac{x^1}{1!}, \frac{x^2}{2!}, \frac{x^3}{3!}, \frac{x^4}{4!}, \cdots$$

against the sequence

$$\cdots, -\frac{3!}{x^4}, \frac{2!}{x^3}, -\frac{1!}{x^2}, \frac{0!}{x^1}, \ln x, \frac{x^1}{1!}, \frac{x^2}{2!}, \frac{x^3}{3!}, \frac{x^4}{4!}, \cdots$$

As $x \rightarrow +\infty$, each sequence increases in magnitude going rightward. Also, each term in each sequence is the derivative with respect to x of the term to its right—except left of the middle element in the first sequence and right of the middle element in the second. The exception is peculiar. What is going on here?

The answer is that x^0 (which is just a constant) and $\ln x$ *both are of zeroth order in x* . This seems strange at first because $\ln x$ diverges as $x \rightarrow \infty$ whereas x^0 does not, but the divergence of the former is extremely slow—so slow, in fact, that per (5.9) $\lim_{x \rightarrow \infty} (\ln x)/x^\epsilon = 0$ for any positive ϵ no matter how small.⁴ Figure 5.2 has plotted $\ln x$ only for $x \sim 1$, but beyond the figure's window the curve (whose slope is $1/x$) flattens rapidly rightward, to the extent that it locally resembles the plot of a constant value; and indeed one can write

$$x^0 = \lim_{u \rightarrow \infty} \frac{\ln(x+u)}{\ln u},$$

which casts x^0 as a logarithm shifted and scaled. Admittedly, one ought not strain such logic too far, because $\ln x$ is not in fact a constant, but the point nevertheless remains that x^0 and $\ln x$ often play analogous roles in mathematics. The logarithm can in some situations profitably be thought of as a “diverging constant” of sorts.

³There are of course some degenerate edge cases like $b = 0$ and $b = 1$. The reader can detail these as the need arises.

⁴One does not grasp how truly slow the divergence is until one calculates a few concrete values. Consider for instance how far out x must run to make $\ln x = 0 \times 100$. It's a long, long way. The natural logarithm does indeed eventually diverge to infinity, in the literal sense that there is no height that it does not eventually reach, but it certainly does not hurry. As we have seen, it takes practically forever just to reach 0×100 .

Less strange-seeming perhaps is the consequence of (5.10) that $\exp x$ is of infinite order in x , that x^∞ and $\exp x$ play analogous roles.

It befits an applied mathematician subjectively to internalize (5.9) and (5.10), to remember that $\ln x$ resembles x^0 and that $\exp x$ resembles x^∞ . A qualitative sense that logarithms are slow and exponentials, fast, helps one to grasp mentally the essential features of many mathematical models one encounters in practice.

Now leaving aside fast and slow functions for the moment, we turn our attention in the next section to the highly important matter of the exponential of a complex argument.

5.4 Euler's formula

The result of § 5.1 leads to one of the central questions in all of mathematics. How can one evaluate

$$\exp i\theta = \lim_{\epsilon \rightarrow 0} (1 + \epsilon)^{i\theta/\epsilon},$$

where $i^2 = -1$ is the imaginary unit introduced in § 2.12?

To begin, one can take advantage of (4.14) to write the last equation in the form

$$\exp i\theta = \lim_{\epsilon \rightarrow 0} (1 + i\epsilon)^{\theta/\epsilon},$$

but from here it is not obvious where to go. The book's development up to the present point gives no obvious direction. In fact it appears that the interpretation of $\exp i\theta$ remains for us to define, if we can find a way to define it which fits sensibly with our existing notions of the real exponential. So, if we don't quite know where to go with this yet, what do we know?

One thing we know is that if $\theta = \epsilon$, then

$$\exp(i\epsilon) = (1 + i\epsilon)^{\epsilon/\epsilon} = 1 + i\epsilon.$$

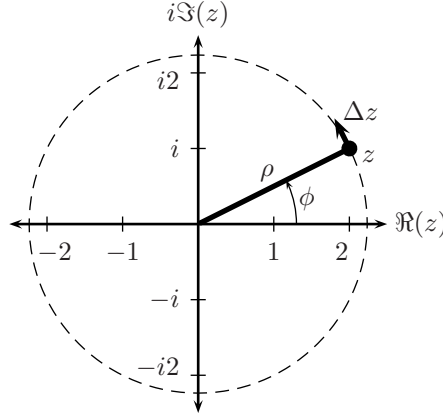
But per § 5.1, the essential operation of the exponential function is to multiply repeatedly by some factor, where the factor is not quite exactly unity—in this case, by $1 + i\epsilon$. So let us multiply a complex number $z = x + iy$ by $1 + i\epsilon$, obtaining

$$(1 + i\epsilon)(x + iy) = (x - \epsilon y) + i(y + \epsilon x).$$

The resulting change in z is

$$\Delta z = (1 + i\epsilon)(x + iy) - (x + iy) = (\epsilon)(-y + ix),$$

Figure 5.3: The complex exponential and Euler's formula.



in which it is seen that

$$\begin{aligned} |\Delta z| &= (\epsilon) \sqrt{y^2 + x^2} = (\epsilon) |z|, \\ \arg(\Delta z) &= \arctan \frac{x}{-y} = \arg z + \frac{2\pi}{4}. \end{aligned}$$

In other words, *the change is proportional to the magnitude of z , but at a right angle to z 's arm in the complex plane.* Refer to Fig. 5.3.

Since the change is directed neither inward nor outward from the dashed circle in the figure, evidently its effect is to move z a distance $(\epsilon) |z|$ counterclockwise about the circle. In other words, referring to the figure, the change Δz leads to

$$\begin{aligned} \Delta \rho &= 0, \\ \Delta \phi &= \epsilon. \end{aligned}$$

These two equations—implying an unvarying, steady counterclockwise progression about the circle in the diagram—are somewhat unexpected. Given the steadiness of the progression, it follows from the equations that

$$\begin{aligned} |\exp i\theta| &= \lim_{\epsilon \rightarrow 0} \left| (1 + i\epsilon)^{\theta/\epsilon} \right| = |\exp 0| + \lim_{\epsilon \rightarrow 0} \frac{\theta}{\epsilon} \Delta \rho = 1, \\ \arg [\exp i\theta] &= \lim_{\epsilon \rightarrow 0} \arg \left[(1 + i\epsilon)^{\theta/\epsilon} \right] = \arg [\exp 0] + \lim_{\epsilon \rightarrow 0} \frac{\theta}{\epsilon} \Delta \phi = \theta. \end{aligned}$$

That is, $\exp i\theta$ is the complex number which lies on the Argand unit circle at phase angle θ . Had we known that θ was an Argand phase angle, naturally we should have represented it by the symbol ϕ from the start. Changing $\phi \leftarrow \theta$ now, we have for real ϕ that

$$\begin{aligned} |\exp i\phi| &= 1, \\ \arg [\exp i\phi] &= \phi, \end{aligned}$$

which says neither more nor less than that

$$\exp i\phi = \cos \phi + i \sin \phi = \text{cis } \phi. \quad (5.11)$$

where $\text{cis}(\cdot)$ is as defined in § 3.11.

Along with the Pythagorean theorem (2.47), the fundamental theorem of calculus (7.2) and Cauchy's integral formula (8.29), (5.11) is one of the most famous results in all of mathematics. It is called *Euler's formula*,^{5,6} and it opens the exponential domain fully to complex numbers, not just for the natural base e but for any base. How? Consider in light of Fig. 5.3 and (5.11) that one can express any complex number in the form

$$z = x + iy = \rho \exp i\phi.$$

If a complex base w is similarly expressed in the form

$$w = u + iv = \sigma \exp i\psi,$$

then it follows that

$$\begin{aligned} w^z &= \exp[\ln w^z] \\ &= \exp[z \ln w] \\ &= \exp[(x + iy)(i\psi + \ln \sigma)] \\ &= \exp[(x \ln \sigma - \psi y) + i(y \ln \sigma + \psi x)]. \end{aligned}$$

Since $\exp(\alpha + \beta) = e^{\alpha+\beta} = \exp \alpha \exp \beta$, the last equation is

$$w^z = \exp(x \ln \sigma - \psi y) \exp i(y \ln \sigma + \psi x), \quad (5.12)$$

⁵For native English speakers who do not speak German, Leonhard Euler's name is pronounced as "oiler."

⁶An alternate derivation of Euler's formula (5.11)—less intuitive and requiring slightly more advanced mathematics, but briefer—constructs from Table 8.1 the Taylor series for $\exp i\phi$, $\cos \phi$ and $i \sin \phi$, then adds the latter two to show them equal to the first of the three. Such an alternate derivation lends little insight, perhaps, but at least it builds confidence that we actually knew what we were doing when we came up with the incredible (5.11).

where

$$\begin{aligned}x &= \rho \cos \phi, \\y &= \rho \sin \phi, \\ \sigma &= \sqrt{u^2 + v^2}, \\ \tan \psi &= \frac{v}{u}.\end{aligned}$$

Equation (5.12) serves to raise any complex number to a complex power.

Curious consequences of Euler's formula (5.11) include that

$$\begin{aligned}e^{\pm i2\pi/4} &= \pm i, \\ e^{\pm i2\pi/2} &= -1, \\ e^{in2\pi} &= 1.\end{aligned}\tag{5.13}$$

For the natural logarithm of a complex number in light of Euler's formula, we have that

$$\ln w = \ln(\sigma e^{i\psi}) = \ln \sigma + i\psi.\tag{5.14}$$

5.5 Complex exponentials and de Moivre's theorem

Euler's formula (5.11) implies that complex numbers z_1 and z_2 can be written

$$\begin{aligned}z_1 &= \rho_1 e^{i\phi_1}, \\ z_2 &= \rho_2 e^{i\phi_2}.\end{aligned}\tag{5.15}$$

By the basic power properties of Table 2.2, then,

$$\begin{aligned}z_1 z_2 &= \rho_1 \rho_2 e^{i(\phi_1 + \phi_2)} = \rho_1 \rho_2 \exp[i(\phi_1 + \phi_2)], \\ \frac{z_1}{z_2} &= \frac{\rho_1}{\rho_2} e^{i(\phi_1 - \phi_2)} = \frac{\rho_1}{\rho_2} \exp[i(\phi_1 - \phi_2)], \\ z^a &= \rho^a e^{ia\phi} = \rho^a \exp[ia\phi].\end{aligned}\tag{5.16}$$

This is de Moivre's theorem, introduced in § 3.11.

5.6 Complex trigonometrics

Applying Euler's formula (5.11) to $+\phi$ then to $-\phi$, we have that

$$\begin{aligned}\exp(+i\phi) &= \cos \phi + i \sin \phi, \\ \exp(-i\phi) &= \cos \phi - i \sin \phi.\end{aligned}$$

Adding the two equations and solving for $\cos \phi$ yields

$$\cos \phi = \frac{\exp(+i\phi) + \exp(-i\phi)}{2}. \quad (5.17)$$

Subtracting the second equation from the first and solving for $\sin \phi$ yields

$$\sin \phi = \frac{\exp(+i\phi) - \exp(-i\phi)}{i2}. \quad (5.18)$$

Thus are the trigonometrics expressed in terms of complex exponentials.

The forms (5.17) and (5.18) suggest the definition of new functions

$$\cosh \phi \equiv \frac{\exp(+\phi) + \exp(-\phi)}{2}, \quad (5.19)$$

$$\sinh \phi \equiv \frac{\exp(+\phi) - \exp(-\phi)}{2}, \quad (5.20)$$

$$\tanh \phi \equiv \frac{\sinh \phi}{\cosh \phi}. \quad (5.21)$$

These are called the *hyperbolic functions*. Their inverses $\operatorname{arccosh}$, etc., are defined in the obvious way. The Pythagorean theorem for trigonometrics (3.2) is that $\cos^2 \phi + \sin^2 \phi = 1$; and from (5.19) and (5.20) one can derive the hyperbolic analog:

$$\begin{aligned} \cos^2 \phi + \sin^2 \phi &= 1, \\ \cosh^2 \phi - \sinh^2 \phi &= 1. \end{aligned} \quad (5.22)$$

The notation $\exp i(\cdot)$ or $e^{i(\cdot)}$ is sometimes felt to be too bulky. Although less commonly seen than the other two, the notation

$$\operatorname{cis}(\cdot) \equiv \exp i(\cdot) = \cos(\cdot) + i \sin(\cdot)$$

is also conventionally recognized, as earlier seen in § 3.11. Also conventionally recognized are $\sin^{-1}(\cdot)$ and occasionally $\operatorname{asin}(\cdot)$ for $\arcsin(\cdot)$, and likewise for the several other trigs.

Replacing $z \leftarrow \phi$ in this section's several equations implies a coherent definition for trigonometric functions of a complex variable. Then, comparing (5.17) and (5.18) respectively to (5.19) and (5.20), we have that

$$\begin{aligned} \cosh z &= \cos iz, \\ i \sinh z &= \sin iz, \\ i \tanh z &= \tan iz, \end{aligned} \quad (5.23)$$

by which one can immediately adapt the many trigonometric properties of Tables 3.1 and 3.3 to hyperbolic use.

At this point in the development one begins to notice that the \sin , \cos , \exp , \csc , \cosh and \sinh functions are each really just different facets of the same mathematical phenomenon. Likewise their respective inverses: \arcsin , \arccos , \ln , $-i \ln$, $\operatorname{arccosh}$ and $\operatorname{arcsinh}$. Conventional names for these two mutually inverse families of functions are unknown to the author, but one might call them the *natural exponential* and *natural logarithmic families*. Or, if the various tangent functions were included, then one might call them the *trigonometric* and *inverse trigonometric families*.

5.7 Summary of properties

Table 5.1 gathers properties of the complex exponential from this chapter and from §§ 2.12, 3.11 and 4.4.

5.8 Derivatives of complex exponentials

This section computes the derivatives of the various trigonometric and inverse trigonometric functions.

5.8.1 Derivatives of sine and cosine

One can compute derivatives of the sine and cosine functions from (5.17) and (5.18), but to do it in that way doesn't seem sporting. Better applied style is to find the derivatives by observing directly the circle from which the sine and cosine functions come.

Refer to Fig. 5.4. Suppose that the point z in the figure is not fixed but travels steadily about the circle such that

$$z(t) = (\rho) [\cos(\omega t + \phi_o) + i \sin(\omega t + \phi_o)]. \quad (5.24)$$

How fast then is the rate dz/dt , and in what Argand direction?

$$\frac{dz}{dt} = (\rho) \left[\frac{d}{dt} \cos(\omega t + \phi_o) + i \frac{d}{dt} \sin(\omega t + \phi_o) \right]. \quad (5.25)$$

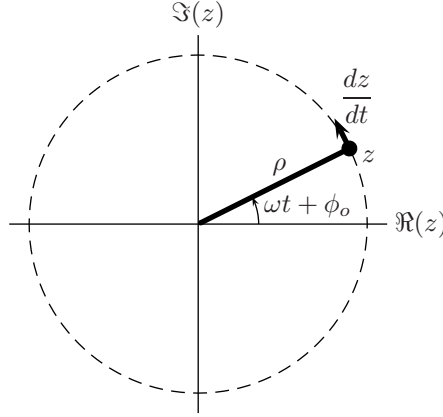
Evidently,

- the speed is $|dz/dt| = (\rho)(d\phi/dt) = \rho\omega$;

Table 5.1: Complex exponential properties.

$$\begin{aligned}
i^2 &= -1 = (-i)^2 \\
\frac{1}{i} &= -i \\
e^{i\phi} &= \cos \phi + i \sin \phi \\
e^{iz} &= \cos z + i \sin z \\
z_1 z_2 &= \rho_1 \rho_2 e^{i(\phi_1 + \phi_2)} = (x_1 x_2 - y_1 y_2) + i(y_1 x_2 + x_1 y_2) \\
\frac{z_1}{z_2} &= \frac{\rho_1}{\rho_2} e^{i(\phi_1 - \phi_2)} = \frac{(x_1 x_2 + y_1 y_2) + i(y_1 x_2 - x_1 y_2)}{x_2^2 + y_2^2} \\
z^a &= \rho^a e^{ia\phi} \\
w^z &= e^{x \ln \sigma - \psi y} e^{i(y \ln \sigma + \psi x)} \\
\ln w &= \ln \sigma + i\psi \\
\sin z &= \frac{e^{iz} - e^{-iz}}{i2} & \sin iz &= i \sinh z & \sinh z &= \frac{e^z - e^{-z}}{2} \\
\cos z &= \frac{e^{iz} + e^{-iz}}{2} & \cos iz &= \cosh z & \cosh z &= \frac{e^z + e^{-z}}{2} \\
\tan z &= \frac{\sin z}{\cos z} & \tan iz &= i \tanh z & \tanh z &= \frac{\sinh z}{\cosh z} \\
\cos^2 z + \sin^2 z &= 1 = \cosh^2 z - \sinh^2 z \\
z &\equiv x + iy = \rho e^{i\phi} \\
w &\equiv u + iv = \sigma e^{i\psi} \\
\exp z &\equiv e^z \\
\operatorname{cis} z &\equiv \cos z + i \sin z = e^{iz} \\
\frac{d}{dz} \exp z &= \exp z \\
\frac{d}{dw} \ln w &= \frac{1}{w} \\
\frac{df/dz}{f(z)} &= \frac{d}{dz} \ln f(z) \\
\log_b w &= \frac{\ln w}{\ln b}
\end{aligned}$$

Figure 5.4: The derivatives of the sine and cosine functions.



- the direction is at right angles to the arm of ρ , which is to say that $\arg(dz/dt) = \phi + 2\pi/4$.

With these observations we can write that

$$\begin{aligned} \frac{dz}{dt} &= (\rho\omega) \left[\cos\left(\omega t + \phi_o + \frac{2\pi}{4}\right) + i \sin\left(\omega t + \phi_o + \frac{2\pi}{4}\right) \right] \\ &= (\rho\omega) [-\sin(\omega t + \phi_o) + i \cos(\omega t + \phi_o)]. \end{aligned} \quad (5.26)$$

Matching the real and imaginary parts of (5.25) against those of (5.26), we have that

$$\begin{aligned} \frac{d}{dt} \cos(\omega t + \phi_o) &= -\omega \sin(\omega t + \phi_o), \\ \frac{d}{dt} \sin(\omega t + \phi_o) &= +\omega \cos(\omega t + \phi_o). \end{aligned} \quad (5.27)$$

If $\omega = 1$ and $\phi_o = 0$, these are

$$\begin{aligned} \frac{d}{dt} \cos t &= -\sin t, \\ \frac{d}{dt} \sin t &= +\cos t. \end{aligned} \quad (5.28)$$

Table 5.2: Derivatives of the trigonometrics.

$$\begin{array}{ll}
\frac{d}{dz} \exp z = + \exp z & \frac{d}{dz} \frac{1}{\exp z} = - \frac{1}{\exp z} \\
\frac{d}{dz} \sin z = + \cos z & \frac{d}{dz} \frac{1}{\sin z} = - \frac{1}{\tan z \sin z} \\
\frac{d}{dz} \cos z = - \sin z & \frac{d}{dz} \frac{1}{\cos z} = + \frac{\tan z}{\cos z}
\end{array}$$

$$\begin{array}{ll}
\frac{d}{dz} \tan z = + (1 + \tan^2 z) & = + \frac{1}{\cos^2 z} \\
\frac{d}{dz} \frac{1}{\tan z} = - \left(1 + \frac{1}{\tan^2 z} \right) & = - \frac{1}{\sin^2 z}
\end{array}$$

$$\begin{array}{ll}
\frac{d}{dz} \sinh z = + \cosh z & \frac{d}{dz} \frac{1}{\sinh z} = - \frac{1}{\tanh z \sinh z} \\
\frac{d}{dz} \cosh z = + \sinh z & \frac{d}{dz} \frac{1}{\cosh z} = - \frac{\tanh z}{\cosh z}
\end{array}$$

$$\begin{array}{ll}
\frac{d}{dz} \tanh z = 1 - \tanh^2 z & = + \frac{1}{\cosh^2 z} \\
\frac{d}{dz} \frac{1}{\tanh z} = 1 - \frac{1}{\tanh^2 z} & = - \frac{1}{\sinh^2 z}
\end{array}$$

5.8.2 Derivatives of the trigonometrics

Equations (5.4) and (5.28) give the derivatives of $\exp(\cdot)$, $\sin(\cdot)$ and $\cos(\cdot)$. From these, with the help of (5.22) and the derivative chain and product rules (§ 4.5), we can calculate the several derivatives of Table 5.2.⁷

⁷[37, back endpaper]

5.8.3 Derivatives of the inverse trigonometrics

Observe the pair

$$\begin{aligned}\frac{d}{dz} \exp z &= \exp z, \\ \frac{d}{dw} \ln w &= \frac{1}{w}.\end{aligned}$$

The natural exponential $\exp z$ belongs to the trigonometric family of functions, as does its derivative. The natural logarithm $\ln w$, by contrast, belongs to the inverse trigonometric family of functions; but its derivative is simpler, not a trigonometric or inverse trigonometric function at all. In Table 5.2, one notices that all the trigonometrics have trigonometric derivatives. By analogy with the natural logarithm, do all the inverse trigonometrics have simpler derivatives?

It turns out that they do. Refer to the account of the natural logarithm's derivative in § 5.2. Following a similar procedure, we have by successive steps that

$$\begin{aligned}\arcsin w &= z, \\ w &= \sin z, \\ \frac{dw}{dz} &= \cos z, \\ \frac{dw}{dz} &= \pm \sqrt{1 - \sin^2 z}, \\ \frac{dw}{dz} &= \pm \sqrt{1 - w^2}, \\ \frac{dz}{dw} &= \frac{\pm 1}{\sqrt{1 - w^2}}, \\ \frac{d}{dw} \arcsin w &= \frac{\pm 1}{\sqrt{1 - w^2}}.\end{aligned}\tag{5.29}$$

Table 5.3: Derivatives of the inverse trigonometrics.

$$\begin{aligned}
\frac{d}{dw} \ln w &= \frac{1}{w} \\
\frac{d}{dw} \arcsin w &= \frac{\pm 1}{\sqrt{1-w^2}} \\
\frac{d}{dw} \arccos w &= \frac{\mp 1}{\sqrt{1-w^2}} \\
\frac{d}{dw} \arctan w &= \frac{1}{1+w^2} \\
\frac{d}{dw} \operatorname{arcsinh} w &= \frac{\pm 1}{\sqrt{w^2+1}} \\
\frac{d}{dw} \operatorname{arccosh} w &= \frac{\pm 1}{\sqrt{w^2-1}} \\
\frac{d}{dw} \operatorname{arctanh} w &= \frac{1}{1-w^2}
\end{aligned}$$

Similarly,

$$\begin{aligned}
\arctan w &= z, \\
w &= \tan z, \\
\frac{dw}{dz} &= 1 + \tan^2 z, \\
\frac{dw}{dz} &= 1 + w^2, \\
\frac{dz}{dw} &= \frac{1}{1+w^2}, \\
\frac{d}{dw} \arctan w &= \frac{1}{1+w^2}. \tag{5.30}
\end{aligned}$$

Derivatives of the other inverse trigonometrics are found in the same way. Table 5.3 summarizes.

5.9 The actuality of complex quantities

Doing all this neat complex math, the applied mathematician can lose sight of some questions he probably ought to keep in mind: Is there really such a

thing as a complex quantity in nature? If not, then hadn't we better avoid these complex quantities, leaving them to the professional mathematical theorists?

As developed by Oliver Heaviside in 1887,⁸ the answer depends on your point of view. If I have 300 g of grapes and 100 g of grapes, then I have 400 g altogether. Alternately, if I have 500 g of grapes and -100 g of grapes, again I have 400 g altogether. (What does it mean to have -100 g of grapes? Maybe that I ate some!) But what if I have $200 + i100$ g of grapes and $200 - i100$ g of grapes? Answer: again, 400 g.

Probably you would not choose to think of $200 + i100$ g of grapes and $200 - i100$ g of grapes, but because of (5.17) and (5.18), one often describes wave phenomena as linear superpositions (sums) of countervailing complex exponentials. Consider for instance the propagating wave

$$A \cos[\omega t - kz] = \frac{A}{2} \exp[+i(\omega t - kz)] + \frac{A}{2} \exp[-i(\omega t - kz)].$$

The benefit of splitting the real cosine into two complex parts is that while the magnitude of the cosine changes with time t , the magnitude of either exponential alone remains steady (see the circle in Fig. 5.3). It turns out to be much easier to analyze two complex wave quantities of constant magnitude than to analyze one real wave quantity of varying magnitude. Better yet, since each complex wave quantity is the complex conjugate of the other, the analyses thereof are mutually conjugate, too; so you normally needn't actually analyze the second. The one analysis suffices for both.⁹ (It's like reflecting your sister's handwriting. To read her handwriting backward, you needn't ask her to try writing reverse with the wrong hand; you can just hold her regular script up to a mirror. Of course, this ignores the question of why one would want to reflect someone's handwriting in the first place; but anyway, reflecting—which is to say, conjugating—complex quantities often is useful.)

Some authors have gently denigrated the use of imaginary parts in physical applications as a mere mathematical trick, as though the parts were not actually there. Well, that is one way to treat the matter, but it is not the way this book recommends. Nothing in the mathematics *requires* you to

⁸[31]

⁹If the point is not immediately clear, an example: Suppose that by the Newton-Raphson iteration (§ 4.8) you have found a root of the polynomial $x^3 + 2x^2 + 3x + 4$ at $x \approx -0.2D + i0.18C$. Where is there another root? Answer: at the complex conjugate, $x \approx -0.2D - i0.18C$. One need not actually run the Newton-Raphson again to find the conjugate root.

regard the imaginary parts as physically nonexistent. You need not abuse Occam's razor! (Occam's razor, "Do not multiply objects without necessity,"¹⁰ is a fine philosophical tool as far as it goes, but is overused in some circles. More often than one likes to believe, the necessity remains hidden until one has ventured to multiply the objects, nor reveals itself to the one who wields the razor, whose hand humility should stay.) It is true by Euler's formula (5.11) that a complex exponential $\exp i\phi$ can be decomposed into a sum of trigonometrics. However, it is equally true by the complex trigonometric formulas (5.17) and (5.18) that *a trigonometric can be decomposed into a sum of complex exponentials*. So, if each can be decomposed into the other, then which of the two is the real decomposition? Answer: it depends on your point of view. Experience seems to recommend viewing the complex exponential as the basic element—as the element of which the trigonometrics are composed—rather than the other way around. From this point of view, it is (5.17) and (5.18) which are the real decomposition. Euler's formula itself is secondary.

The complex exponential method of offsetting imaginary parts offers an elegant yet practical mathematical way to model physical wave phenomena. So go ahead: regard the imaginary parts as actual. It doesn't hurt anything, and it helps with the math.

¹⁰[42, Ch. 12]

Chapter 6

Primes, roots and averages

This chapter gathers a few significant topics, each of whose treatment seems too brief for a chapter of its own.

6.1 Prime numbers

A *prime number*—or simply, a *prime*—is an integer greater than one, divisible only by one and itself. A *composite number* is an integer greater than one and not prime. A composite number can be composed as a product of two or more prime numbers. All positive integers greater than one are either composite or prime.

The mathematical study of prime numbers and their incidents constitutes *number theory*, and it is a deep area of mathematics. The deeper results of number theory seldom arise in applications,¹ however, so we shall confine our study of number theory in this book to one or two of its simplest, most broadly interesting results.

6.1.1 The infinite supply of primes

The first primes are evidently 2, 3, 5, 7, 11, Is there a last prime? To show that there is not, suppose that there were. More precisely, suppose that there existed exactly N primes, with N finite, letting p_1, p_2, \dots, p_N represent these primes from least to greatest. Now consider the product of

¹The deeper results of number theory do arise in cryptography, or so the author has been led to understand. Although cryptography is literally an application of mathematics, its spirit is that of pure mathematics rather than of applied. If you seek cryptographic derivations, this book is probably not the one you want.

all the primes,

$$C = \prod_{j=1}^N p_j.$$

What of $C + 1$? Since $p_1 = 2$ divides C , it cannot divide $C + 1$. Similarly, since $p_2 = 3$ divides C , it also cannot divide $C + 1$. The same goes for $p_3 = 5$, $p_4 = 7$, $p_5 = 11$, etc. Apparently none of the primes in the p_j series divides $C + 1$, which implies either that $C + 1$ itself is prime, or that $C + 1$ is composed of primes not in the series. But the latter is assumed impossible on the ground that the p_j series includes all primes; and the former is assumed impossible on the ground that $C + 1 > C > p_N$, with p_N the greatest prime. The contradiction proves false the assumption which gave rise to it. The false assumption: that there were a last prime.

Thus there is no last prime. No matter how great a prime number one finds, a greater can always be found. The supply of primes is infinite.²

Attributed to the ancient geometer Euclid, the foregoing proof is a classic example of mathematical *reductio ad absurdum*, or as usually styled in English, *proof by contradiction*.³

6.1.2 Compositional uniqueness

Occasionally in mathematics, plausible assumptions can hide subtle logical flaws. One such plausible assumption is the assumption that every positive integer has a unique *prime factorization*. It is readily seen that the first several positive integers— $1 = ()$, $2 = (2^1)$, $3 = (3^1)$, $4 = (2^2)$, $5 = (5^1)$, $6 = (2^1)(3^1)$, $7 = (7^1)$, $8 = (2^3)$, \dots —each have unique prime factorizations, but is this necessarily true of all positive integers?

To show that it is true, suppose that it were not.⁴ More precisely, suppose that there did exist positive integers factorable each in two or more distinct ways, with the symbol C representing the least such integer. Noting that C must be composite (prime numbers by definition are each factorable

²[40]

³[34, Appendix 1][46, “Reductio ad absurdum,” 02:36, 28 April 2006]

⁴Unfortunately the author knows no more elegant proof than this, yet cannot even cite this one properly. The author encountered the proof in some book over a decade ago. The identity of that book is now long forgotten.

only one way, like $5 = [5^1]$), let

$$\begin{aligned}
 C_p &\equiv \prod_{j=1}^{N_p} p_j, \\
 C_q &\equiv \prod_{k=1}^{N_q} q_k, \\
 C_p = C_q &= C, \\
 p_j &\leq p_{j+1}, \\
 q_k &\leq q_{k+1}, \\
 p_1 &\leq q_1, \\
 N_p &> 1, \\
 N_q &> 1,
 \end{aligned}$$

where C_p and C_q represent two distinct prime factorizations of the same number C and where the p_j and q_k are the respective primes ordered from least to greatest. We see that

$$p_j \neq q_k$$

for any j and k —that is, that the same prime cannot appear in both factorizations—because if the same prime r did appear in both then C/r either would be prime (in which case both factorizations would be $[r][C/r]$, defying our assumption that the two differed) or would constitute an ambiguously factorable composite integer less than C when we had already defined C to represent the least such. Among other effects, the fact that $p_j \neq q_k$ strengthens the definition $p_1 \leq q_1$ to read

$$p_1 < q_1.$$

Let us now rewrite the two factorizations in the form

$$\begin{aligned}
 C_p &= p_1 A_p, \\
 C_q &= q_1 A_q, \\
 C_p = C_q &= C, \\
 A_p &\equiv \prod_{j=2}^{N_p} p_j, \\
 A_q &\equiv \prod_{k=2}^{N_q} q_k,
 \end{aligned}$$

where p_1 and q_1 are the least primes in their respective factorizations. Since C is composite and since $p_1 < q_1$, we have that

$$1 < p_1 < q_1 \leq \sqrt{C} \leq A_q < A_p < C,$$

which implies that

$$p_1 q_1 < C.$$

The last inequality lets us compose the new positive integer

$$B = C - p_1 q_1,$$

which might be prime or composite (or unity), but which either way enjoys a unique prime factorization because $B < C$, with C the least positive integer factorable two ways. Observing that some integer s which divides C necessarily also divides $C \pm ns$, we note that each of p_1 and q_1 necessarily divides B . This means that B 's unique factorization includes both p_1 and q_1 , which further means that the product $p_1 q_1$ divides B . But if $p_1 q_1$ divides B , then it divides $B + p_1 q_1 = C$, also.

Let E represent the positive integer which results from dividing C by $p_1 q_1$:

$$E \equiv \frac{C}{p_1 q_1}.$$

Then,

$$\begin{aligned} E q_1 &= \frac{C}{p_1} = A_p, \\ E p_1 &= \frac{C}{q_1} = A_q. \end{aligned}$$

That $E q_1 = A_p$ says that q_1 divides A_p . But $A_p < C$, so A_p 's prime factorization is unique—and we see above that A_p 's factorization *does not include any* q_k , not even q_1 . The contradiction proves false the assumption which gave rise to it. The false assumption: that there existed a least composite number C prime-factorable in two distinct ways.

Thus no positive integer is ambiguously factorable. Prime factorizations are always unique.

We have observed at the start of this subsection that plausible assumptions can hide subtle logical flaws. Indeed this is so. Interestingly however, the plausible assumption of the present subsection has turned out absolutely correct; we have just had to do some extra work to prove it. Such effects

are typical on the shadowed frontier where applied shades into pure mathematics: with sufficient experience and with a firm grasp of the model at hand, if you think that it's true, then it probably is. Judging when to delve into the mathematics anyway, seeking a more rigorous demonstration of a proposition one feels pretty sure is correct, is a matter of applied mathematical style. It depends on how sure one feels, and more importantly on whether the unsureness felt is true uncertainty or is just an unaccountable desire for more precise mathematical definition (if the latter, then unlike the author you may have the right temperament to become a professional mathematician). The author does judge the present subsection's proof to be worth the applied effort; but nevertheless, when one lets logical minutiae distract him to too great a degree, he admittedly begins to drift out of the applied mathematical realm that is the subject of this book.

6.1.3 Rational and irrational numbers

A *rational number* is a finite real number expressible as a ratio of integers

$$x = \frac{p}{q}, \quad q \neq 0.$$

The ratio is *fully reduced* if p and q have no prime factors in common. For instance, $4/6$ is not fully reduced, whereas $2/3$ is.

An *irrational number* is a finite real number which is not rational. For example, $\sqrt{2}$ is irrational. In fact any $x = \sqrt{n}$ is irrational unless integral; there is no such thing as a \sqrt{n} which is not an integer but is rational.

To prove⁵ the last point, suppose that there did exist a fully reduced

$$x = \frac{p}{q} = \sqrt{n}, \quad p > 0, \quad q > 1,$$

where p , q and n are all integers. Squaring the equation, we have that

$$\frac{p^2}{q^2} = n,$$

which form is evidently also fully reduced. But if $q > 1$, then the fully reduced $n = p^2/q^2$ is not an integer as we had assumed that it was. The contradiction proves false the assumption which gave rise to it. Hence there exists no rational, nonintegral \sqrt{n} , as was to be demonstrated. The proof is readily extended to show that any $x = n^{j/k}$ is irrational if nonintegral, the

⁵A proof somewhat like the one presented here is found in [34, Appendix 1].

extension by writing $p^k/q^k = n^j$ then following similar steps as those this paragraph outlines.

That's all the number theory the book treats; but in applied math, so little will take you pretty far. Now onward we go to other topics.

6.2 The existence and number of polynomial roots

This section shows that an N th-order polynomial must have exactly N roots.

6.2.1 Polynomial roots

Consider the quotient $B(z)/A(z)$, where

$$\begin{aligned} A(z) &= z - \alpha, \\ B(z) &= \sum_{k=0}^N b_k z^k, \quad N > 0, \quad b_N \neq 0, \\ B(\alpha) &= 0. \end{aligned}$$

In the long-division symbology of Table 2.3,

$$B(z) = A(z)Q_0(z) + R_0(z),$$

where $Q_0(z)$ is the quotient and $R_0(z)$, a remainder. In this case the divisor $A(z) = z - \alpha$ has first order, and as § 2.6.2 has observed, first-order divisors leave zeroth-order, constant remainders $R_0(z) = \rho$. Thus substituting yields

$$B(z) = (z - \alpha)Q_0(z) + \rho.$$

When $z = \alpha$, this reduces to

$$B(\alpha) = \rho.$$

But $B(\alpha) = 0$, so

$$\rho = 0.$$

Evidently the division leaves no remainder ρ , which is to say that $z - \alpha$ *exactly divides every polynomial $B(z)$ of which $z = \alpha$ is a root.*

Note that if the polynomial $B(z)$ has order N , then the quotient $Q(z) = B(z)/(z - \alpha)$ has exactly order $N - 1$. That is, the leading, z^{N-1} term of the quotient is never null. The reason is that if the leading term were null, if $Q(z)$ had order less than $N - 1$, then $B(z) = (z - \alpha)Q(z)$ could not possibly have order N as we have assumed.

6.2.2 The fundamental theorem of algebra

The *fundamental theorem of algebra* holds that any polynomial $B(z)$ of order N can be factored

$$B(z) = \sum_{k=0}^N b_k z^k = b_N \prod_{j=1}^N (z - \alpha_j), \quad b_N \neq 0, \quad (6.1)$$

where the α_k are the N roots of the polynomial.⁶

To prove the theorem, it suffices to show that all polynomials of order $N > 0$ have at least one root; for if a polynomial of order N has a root α_N , then according to § 6.2.1 one can divide the polynomial by $z - \alpha_N$ to obtain a new polynomial of order $N - 1$. To the new polynomial the same logic applies; if it has at least one root α_{N-1} , then one can divide it by $z - \alpha_{N-1}$ to obtain yet another polynomial of order $N - 2$; and so on, one root extracted at each step, factoring the polynomial step by step into the desired form $b_N \prod_{j=1}^N (z - \alpha_j)$.

It remains however to show that there exists no polynomial $B(z)$ of order $N > 0$ lacking roots altogether. To show that there is no such polynomial, consider the locus⁷ of all $B(\rho e^{i\phi})$ in the Argand range plane (Fig. 2.5), where $z = \rho e^{i\phi}$, ρ is held constant, and ϕ is variable. Because $e^{i(\phi+n2\pi)} = e^{i\phi}$ and no fractional powers of z appear in (6.1), this locus forms a closed loop. At very large ρ , the $b_N z^N$ term dominates $B(z)$, so the locus there evidently has the general character of $b_N \rho^N e^{iN\phi}$. As such, the locus is nearly but not quite a circle at radius $b_N \rho^N$ from the Argand origin $B(z) = 0$, revolving N times at that great distance before exactly repeating. On the other hand, when $\rho = 0$ the entire locus collapses on the single point $B(0) = b_0$.

Now consider the locus at very large ρ again, but this time let ρ slowly shrink. Watch the locus as ρ shrinks. The locus is like a great string or rubber band, joined at the ends and looped in N great loops. As ρ shrinks smoothly, the string's shape changes smoothly. Eventually ρ disappears and the entire string collapses on the point $B(0) = b_0$. Since the string originally has looped N times at great distance about the Argand origin, but at the end has collapsed on a single point, then at some time between it must have swept through the origin and every other point within the original loops.

⁶Professional mathematicians typically state the theorem in a slightly different form. They also prove it in rather a different way. [21, Ch. 10, Prob. 74]

⁷A *locus* is the geometric collection of points which satisfy a given criterion. For example, the locus of all points in a plane at distance ρ from a point O is a circle; the locus of all points in three-dimensional space equidistant from two points P and Q is a plane; etc.

After all, $B(z)$ is everywhere differentiable, so the string can only *sweep* as ρ decreases; it can never skip. The Argand origin lies inside the loops at the start but outside at the end. If so, then the values of ρ and ϕ precisely where the string has swept through the origin by definition constitute a root $B(\rho e^{i\phi}) = 0$. Thus as we were required to show, $B(z)$ does have at least one root, which observation completes the applied demonstration of the fundamental theorem of algebra.

The fact that the roots exist is one thing. Actually finding the roots numerically is another matter. For a quadratic (second order) polynomial, (2.2) gives the roots. For cubic (third order) and quartic (fourth order) polynomials, formulas for the roots are known (see Ch. 10) though seemingly not so for quintic (fifth order) and higher-order polynomials;⁸ but the Newton-Raphson iteration (§ 4.8) can be used to locate a root numerically in any case. The Newton-Raphson is used to extract one root (*any* root) at each step as described above, reducing the polynomial step by step until all the roots are found.

The reverse problem, finding the polynomial given the roots, is much easier: one just multiplies out $\prod_j (z - \alpha_j)$, as in (6.1).

6.3 Addition and averages

This section discusses the two basic ways to add numbers and the three basic ways to calculate averages of them.

6.3.1 Serial and parallel addition

Consider the following problem. There are three masons. The strongest and most experienced of the three, Adam, lays 120 bricks per hour.⁹ Next is Brian who lays 90. Charles is new; he lays only 60. Given eight hours, how many bricks can the three men lay? Answer:

$$(8 \text{ hours})(120 + 90 + 60 \text{ bricks per hour}) = 2160 \text{ bricks.}$$

Now suppose that we are told that Adam can lay a brick every 30 seconds; Brian, every 40 seconds; Charles, every 60 seconds. How much time do the

⁸In a celebrated theorem of pure mathematics [44, “Abel’s impossibility theorem”], it is said to be shown that no such formula even exists, given that the formula be constructed according to certain rules. Undoubtedly the theorem is interesting to the professional mathematician, but to the applied mathematician it probably suffices to observe merely that no such formula is known.

⁹The figures in the example are in decimal notation.

three men need to lay 2160 bricks? Answer:

$$\begin{aligned} \frac{2160 \text{ bricks}}{\frac{1}{30} + \frac{1}{40} + \frac{1}{60} \text{ bricks per second}} &= 28,800 \text{ seconds} \left(\frac{1 \text{ hour}}{3600 \text{ seconds}} \right) \\ &= 8 \text{ hours.} \end{aligned}$$

The two problems are precisely equivalent. Neither is stated in simpler terms than the other. The notation used to solve the second is less elegant, but fortunately there exists a better notation:

$$(2160 \text{ bricks})(30 \parallel 40 \parallel 60 \text{ seconds per brick}) = 8 \text{ hours,}$$

where

$$\frac{1}{30 \parallel 40 \parallel 60} = \frac{1}{30} + \frac{1}{40} + \frac{1}{60}.$$

The operator \parallel is called the *parallel addition* operator. It works according to the law

$$\frac{1}{a \parallel b} = \frac{1}{a} + \frac{1}{b}, \quad (6.2)$$

where the familiar operator $+$ is verbally distinguished from the \parallel when necessary by calling the $+$ the *serial addition* or *series addition* operator. With (6.2) and a bit of arithmetic, the several parallel-addition identities of Table 6.1 are soon derived.

The writer knows of no conventional notation for parallel sums of series, but suggests that the notation which appears in the table,

$$\sum_{k=a}^b \parallel f(k) \equiv f(a) \parallel f(a+1) \parallel f(a+2) \parallel \cdots \parallel f(b),$$

might serve if needed.

Assuming that none of the values involved is negative, one can readily show that¹⁰

$$a \parallel x \leq b \parallel x \text{ iff } a \leq b. \quad (6.3)$$

This is intuitive. Counterintuitive, perhaps, is that

$$a \parallel x \leq a. \quad (6.4)$$

Because we have all learned as children to count in the sensible manner 1, 2, 3, 4, 5, ...—rather than as 1, $\frac{1}{2}$, $\frac{1}{3}$, $\frac{1}{4}$, $\frac{1}{5}$, ...—serial addition ($+$) seems

¹⁰The word *iff* means, “if and only if.”

Table 6.1: Parallel and serial addition identities.

$\frac{1}{a \parallel b} = \frac{1}{a} + \frac{1}{b}$	$\frac{1}{a + b} = \frac{1}{a} \parallel \frac{1}{b}$
$a \parallel b = \frac{ab}{a + b}$	$a + b = \frac{ab}{a \parallel b}$
$a \parallel \frac{1}{b} = \frac{a}{1 + ab}$	$a + \frac{1}{b} = \frac{a}{1 \parallel ab}$
$a \parallel b = b \parallel a$	$a + b = b + a$
$a \parallel (b \parallel c) = (a \parallel b) \parallel c$	$a + (b + c) = (a + b) + c$
$a \parallel \infty = \infty \parallel a = a$	$a + 0 = 0 + a = a$
$a \parallel (-a) = \infty$	$a + (-a) = 0$
$(a)(b \parallel c) = ab \parallel ac$	$(a)(b + c) = ab + ac$
$\frac{1}{\sum_k \parallel a_k} = \sum_k \frac{1}{a_k}$	$\frac{1}{\sum_k a_k} = \sum_k \parallel \frac{1}{a_k}$

more natural than parallel addition (\parallel) does. The psychological barrier is hard to breach, yet for many purposes parallel addition is in fact no less fundamental. Its rules are inherently neither more nor less complicated, as Table 6.1 illustrates; yet outside the electrical engineering literature the parallel addition notation is seldom seen.¹¹ Now that you have seen it, you can use it. There is profit in learning to think both ways. (Exercise: counting from zero serially goes 0, 1, 2, 3, 4, 5, ...; how does the parallel analog go?)¹²

Convention brings no special symbol for parallel subtraction, incidentally. One merely writes

$$a \parallel (-b),$$

which means exactly what it appears to mean.

6.3.2 Averages

Let us return to the problem of the preceding section. Among the three masons, what is their average productivity? The answer depends on how you look at it. On the one hand,

$$\frac{120 + 90 + 60 \text{ bricks per hour}}{3} = 90 \text{ bricks per hour.}$$

On the other hand,

$$\frac{30 + 40 + 60 \text{ seconds per brick}}{3} = 43\frac{1}{3} \text{ seconds per brick.}$$

These two figures are not the same. That is, $1/(43\frac{1}{3} \text{ seconds per brick}) \neq 90 \text{ bricks per hour}$. Yet both figures are valid. Which figure you choose depends on what you want to calculate. A common mathematical error among businesspeople is not to realize that both averages are possible and that they yield different numbers (if the businessperson quotes in bricks per hour, the productivities average one way; if in seconds per brick, the other way; yet some businesspeople will never clearly consider the difference). Realizing this, the clever businessperson might negotiate a contract so that the average used worked to his own advantage.¹³

¹¹In electric circuits, loads are connected in parallel as often as, in fact probably more often than, they are connected in series. Parallel addition gives the electrical engineer a neat way of adding the impedances of parallel-connected loads.

¹²[36, eqn. 1.27]

¹³“And what does the author know about business?” comes the rejoinder.

The rejoinder is fair enough. If the author wanted to demonstrate his business acumen (or lack thereof) he’d do so elsewhere not here! There are a lot of good business books

When it is unclear which of the two averages is more appropriate, a third average is available, the *geometric mean*

$$[(120)(90)(60)]^{1/3} \text{ bricks per hour.}$$

The geometric mean does not have the problem either of the two averages discussed above has. The inverse geometric mean

$$[(30)(40)(60)]^{1/3} \text{ seconds per brick}$$

implies the same average productivity. The mathematically savvy sometimes prefer the geometric mean over either of the others for this reason.

Generally, the *arithmetic*, *geometric* and *harmonic means* are defined

$$\mu \equiv \frac{\sum_k w_k x_k}{\sum_k w_k} = \left(\sum_k \parallel \frac{1}{w_k} \right) \left(\sum_k w_k x_k \right), \quad (6.5)$$

$$\mu_{\Pi} \equiv \left[\prod_j x_j^{w_j} \right]^{1/\sum_k w_k} = \left[\prod_j x_j^{w_j} \right]^{\sum_k \parallel 1/w_k}, \quad (6.6)$$

$$\mu_{\parallel} \equiv \frac{\sum_k \parallel x_k/w_k}{\sum_k \parallel 1/w_k} = \left(\sum_k w_k \right) \left(\sum_k \parallel \frac{x_k}{w_k} \right), \quad (6.7)$$

where the x_k are the several samples and the w_k are weights. For two samples weighted equally, these are

$$\mu = \frac{a+b}{2}, \quad (6.8)$$

$$\mu_{\Pi} = \sqrt{ab}, \quad (6.9)$$

$$\mu_{\parallel} = 2(a \parallel b). \quad (6.10)$$

out there and this is not one of them.

The fact remains nevertheless that businesspeople sometimes use mathematics in peculiar ways, making relatively easy problems harder and more mysterious than the problems need to be. If you have ever encountered the little monstrosity of an approximation banks (at least in the author's country) actually use in place of (9.12) to accrue interest and amortize loans, then you have met the difficulty.

Trying to convince businesspeople that their math is wrong, incidentally, is in the author's experience usually a waste of time. Some businesspeople are mathematically rather sharp—as you presumably are if you are in business and are reading these words—but as for most: when real mathematical ability is needed, that's what they hire engineers, architects and the like for. The author is not sure, but somehow he doubts that many boards of directors would be willing to bet the company on a financial formula containing some mysterious-looking e^x . Business demands other talents.

If $a \geq 0$ and $b \geq 0$, then by successive steps,¹⁴

$$\begin{aligned}
 0 &\leq (a - b)^2, \\
 0 &\leq a^2 - 2ab + b^2, \\
 4ab &\leq a^2 + 2ab + b^2, \\
 2\sqrt{ab} &\leq a + b, \\
 \frac{2\sqrt{ab}}{a + b} &\leq 1 \leq \frac{a + b}{2\sqrt{ab}}, \\
 \frac{2ab}{a + b} &\leq \sqrt{ab} \leq \frac{a + b}{2}, \\
 2(a \parallel b) &\leq \sqrt{ab} \leq \frac{a + b}{2}.
 \end{aligned}$$

That is,

$$\mu_{\parallel} \leq \mu_{\Pi} \leq \mu. \quad (6.11)$$

The arithmetic mean is greatest and the harmonic mean, least; with the geometric mean falling between.

Does (6.11) hold when there are several nonnegative samples of various nonnegative weights? To show that it does, consider the case of $N = 2^m$ nonnegative samples of equal weight. Nothing prevents one from dividing such a set of samples in half, considering each subset separately, for if (6.11) holds for each subset individually then surely it holds for the whole set (this is so because the average of the whole set is itself the *average of the two subset averages*, where the word “average” signifies the arithmetic, geometric or harmonic mean as appropriate). But each subset can further be divided in half, then each subsubset can be divided in half again, and so on until each smallest group has two members only—in which case we already know that (6.11) obtains. Starting there and recursing back, we have that (6.11)

¹⁴The steps are logical enough, but the motivation behind them remains inscrutable until the reader realizes that the writer originally worked the steps out backward with his pencil, from the last step to the first. Only then did he reverse the order and write the steps formally here. The writer had no idea that he was supposed to start from $0 \leq (a - b)^2$ until his pencil working backward showed him. “Begin with the end in mind,” the saying goes. In this case the saying is right.

The same reading strategy often clarifies inscrutable math. When you can follow the logic but cannot understand what could possibly have inspired the writer to conceive the logic in the first place, try reading backward.

obtains for the entire set. Now consider that a sample of any weight can be approximated arbitrarily closely by several samples of weight $1/2^m$, provided that m is sufficiently large. By this reasoning, (6.11) holds for any nonnegative weights of nonnegative samples, which was to be demonstrated.

Chapter 7

The integral

Chapter 4 has observed that the mathematics of calculus concerns a complementary pair of questions:

- Given some function $f(t)$, what is the function's instantaneous rate of change, or *derivative*, $f'(t)$?
- Interpreting some function $f'(t)$ as an instantaneous rate of change, what is the corresponding accretion, or *integral*, $f(t)$?

Chapter 4 has built toward a basic understanding of the first question. This chapter builds toward a basic understanding of the second. The understanding of the second question constitutes the concept of the integral, one of the profoundest ideas in all of mathematics.

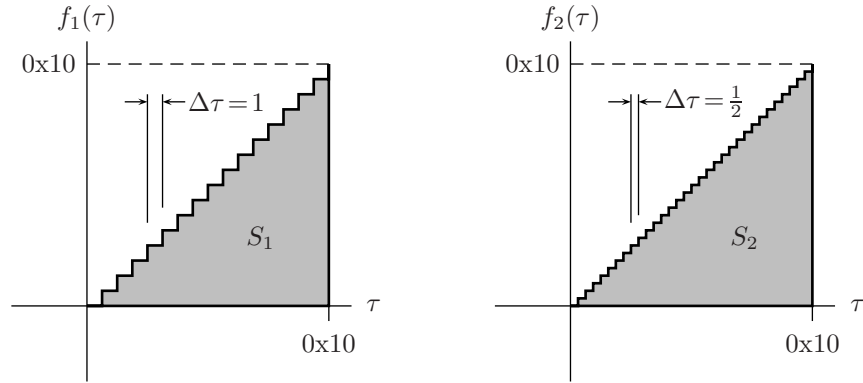
This chapter, which introduces the integral, is undeniably a hard chapter.

Experience knows no reliable way to teach the integral adequately to the uninitiated except through dozens or hundreds of pages of suitable examples and exercises, yet the book you are reading cannot be that kind of book. The sections of the present chapter concisely treat matters which elsewhere rightly command chapters or whole books of their own. Concision can be a virtue—and by design, nothing essential is omitted here—but the bold novice who wishes to learn the integral from these pages alone faces a daunting challenge. It can be done. However, for less intrepid readers who quite reasonably prefer a gentler initiation, [18] is warmly recommended.

7.1 The concept of the integral

An *integral* is a finite accretion or sum of an infinite number of infinitesimal elements. This section introduces the concept.

Figure 7.1: Areas representing discrete sums.



7.1.1 An introductory example

Consider the sums

$$\begin{aligned}
 S_1 &= \sum_{k=0}^{0x10-1} k, \\
 S_2 &= \frac{1}{2} \sum_{k=0}^{0x20-1} \frac{k}{2}, \\
 S_4 &= \frac{1}{4} \sum_{k=0}^{0x40-1} \frac{k}{4}, \\
 S_8 &= \frac{1}{8} \sum_{k=0}^{0x80-1} \frac{k}{8}, \\
 &\vdots \\
 S_n &= \frac{1}{n} \sum_{k=0}^{(0x10)n-1} \frac{k}{n}.
 \end{aligned}$$

What do these sums represent? One way to think of them is in terms of the shaded areas of Fig. 7.1. In the figure, S_1 is composed of several tall, thin rectangles of width 1 and height k ; S_2 , of rectangles of width $1/2$ and height

$k/2$.¹ As n grows, the shaded region in the figure looks more and more like a triangle of base length $b = 0x10$ and height $h = 0x10$. In fact it appears that

$$\lim_{n \rightarrow \infty} S_n = \frac{bh}{2} = 0x80,$$

or more tersely

$$S_\infty = 0x80,$$

is the area the increasingly fine stairsteps approach.

Notice how we have evaluated S_∞ , the sum of an infinite number of infinitely narrow rectangles, without actually adding anything up. We have taken a shortcut directly to the total.

In the equation

$$S_n = \frac{1}{n} \sum_{k=0}^{(0x10)n-1} \frac{k}{n},$$

let us now change the variables

$$\begin{aligned} \tau &\leftarrow \frac{k}{n}, \\ \Delta\tau &\leftarrow \frac{1}{n}, \end{aligned}$$

to obtain the representation

$$S_n = \Delta\tau \sum_{k=0}^{(0x10)n-1} \tau;$$

or more properly,

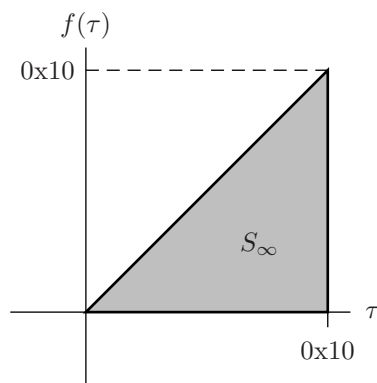
$$S_n = \sum_{k=0}^{(k|_{\tau=0x10})-1} \tau \Delta\tau,$$

where the notation $k|_{\tau=0x10}$ indicates the value of k when $\tau = 0x10$. Then

$$S_\infty = \lim_{\Delta\tau \rightarrow 0^+} \sum_{k=0}^{(k|_{\tau=0x10})-1} \tau \Delta\tau,$$

¹If the reader does not fully understand this paragraph's illustration, if the relation of the sum to the area seems unclear, the reader is urged to pause and consider the illustration carefully until he does understand it. If it still seems unclear, then the reader should probably suspend reading here and go study a good basic calculus text like [18]. The concept is important.

Figure 7.2: An area representing an infinite sum of infinitesimals. (Observe that the infinitesimal $d\tau$ is now too narrow to show on this scale. Compare against $\Delta\tau$ in Fig. 7.1.)



in which it is conventional as $\Delta\tau$ vanishes to change the symbol $d\tau \leftarrow \Delta\tau$, where $d\tau$ is the infinitesimal of Ch. 4:

$$S_\infty = \lim_{d\tau \rightarrow 0^+} \sum_{k=0}^{(k|_{\tau=0x10})-1} \tau d\tau.$$

The symbol $\lim_{d\tau \rightarrow 0^+} \sum_{k=0}^{(k|_{\tau=0x10})-1}$ is cumbersome, so we replace it with the new symbol² \int_0^{0x10} to obtain the form

$$S_\infty = \int_0^{0x10} \tau d\tau.$$

This means, “stepping in infinitesimal intervals of $d\tau$, the sum of all $\tau d\tau$ from $\tau = 0$ to $\tau = 0x10$.” Graphically, it is the shaded area of Fig. 7.2.

²Like the Greek S, \sum , denoting discrete summation, the seventeenth century-styled Roman S, \int , stands for Latin “summa,” English “sum.” See [46, “Long s,” 14:54, 7 April 2006].

7.1.2 Generalizing the introductory example

Now consider a generalization of the example of § 7.1.1:

$$S_n = \frac{1}{n} \sum_{k=an}^{bn-1} f\left(\frac{k}{n}\right).$$

(In the example of § 7.1.1, $f[\tau]$ was the simple $f[\tau] = \tau$, but in general it could be any function.) With the change of variables

$$\begin{aligned}\tau &\leftarrow \frac{k}{n}, \\ \Delta\tau &\leftarrow \frac{1}{n},\end{aligned}$$

this is

$$S_n = \sum_{k=(k|_{\tau=a})}^{(k|_{\tau=b})-1} f(\tau) \Delta\tau.$$

In the limit,

$$S_\infty = \lim_{d\tau \rightarrow 0^+} \sum_{k=(k|_{\tau=a})}^{(k|_{\tau=b})-1} f(\tau) d\tau = \int_a^b f(\tau) d\tau.$$

This is the *integral* of $f(\tau)$ in the interval $a < \tau < b$. It represents the area under the curve of $f(\tau)$ in that interval.

7.1.3 The balanced definition and the trapezoid rule

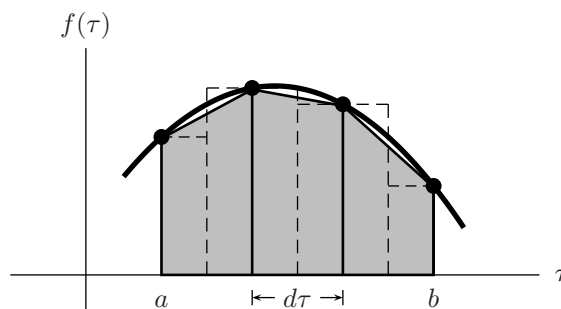
Actually, just as we have defined the derivative in the balanced form (4.15), we do well to define the integral in balanced form, too:

$$\int_a^b f(\tau) d\tau \equiv \lim_{d\tau \rightarrow 0^+} \left\{ \frac{f(a) d\tau}{2} + \sum_{k=(k|_{\tau=a})+1}^{(k|_{\tau=b})-1} f(\tau) d\tau + \frac{f(b) d\tau}{2} \right\}. \quad (7.1)$$

Here, the first and last integration samples are each balanced “on the edge,” half within the integration domain and half without.

Equation (7.1) is known as the *trapezoid rule*. Figure 7.3 depicts it. The name “trapezoid” comes of the shapes of the shaded integration elements in the figure. Observe however that it makes no difference whether one regards the shaded trapezoids or the dashed rectangles as the actual integration

Figure 7.3: Integration by the trapezoid rule (7.1). Notice that the shaded and dashed areas total the same.



elements; the total integration area is the same either way.³ The important point to understand is that the integral is conceptually just a sum. It is a sum of an infinite number of infinitesimal elements as $d\tau$ tends to vanish, but a sum nevertheless; nothing more.

Nothing actually requires the integration element width $d\tau$ to remain constant from element to element, incidentally. Constant widths are usually easiest to handle but variable widths find use in some cases. The only requirement is that $d\tau$ remain infinitesimal. (For further discussion of the point, refer to the treatment of the Leibnitz notation in § 4.4.2.)

7.2 The antiderivative and the fundamental theorem of calculus

If

$$S(x) \equiv \int_a^x g(\tau) d\tau,$$

³The trapezoid rule (7.1) is perhaps the most straightforward, general, robust way to define the integral, but other schemes are possible, too. For example, one can give the integration elements quadratically curved tops which more nearly track the actual curve. That scheme is called *Simpson's rule*. [A section on Simpson's rule might be added to the book at some later date.]

then what is the derivative dS/dx ? After some reflection, one sees that the derivative must be

$$\frac{dS}{dx} = g(x).$$

This is so because the action of the integral is to compile or accrete the area under a curve. The integral accretes area at a rate proportional to the curve's height $f(\tau)$: the higher the curve, the faster the accretion. In this way one sees that the integral and the derivative are inverse operators; the one inverts the other. The integral is the *antiderivative*.

More precisely,

$$\int_a^b \frac{df}{d\tau} d\tau = f(\tau)|_a^b, \quad (7.2)$$

where the notation $f(\tau)|_a^b$ or $[f(\tau)]_a^b$ means $f(b) - f(a)$.

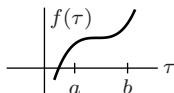
The importance of (7.2), fittingly named the *fundamental theorem of calculus*,⁴ can hardly be overstated. As the formula which ties together the complementary pair of questions asked at the chapter's start, (7.2) is of utmost importance in the practice of mathematics. The idea behind the formula is indeed simple once grasped, but to grasp the idea firmly in the first place is not entirely trivial.⁵ The idea is simple but big. The reader

⁴[18, § 11.6][37, § 5-4][46, "Fundamental theorem of calculus," 06:29, 23 May 2006]

⁵Having read from several calculus books and, like millions of others perhaps including the reader, having sat years ago in various renditions of the introductory calculus lectures in school, the author has never yet met a more convincing demonstration of (7.2) than the formula itself. Somehow the underlying idea is too simple, too profound to explain. It's like trying to explain how to drink water, or how to count or to add. Elaborate explanations and their attendant constructs and formalities are indeed possible to contrive, but the idea itself is so simple that somehow such contrivances seem to obscure the idea more than to reveal it.

One ponders the formula (7.2) a while, then the idea dawns on him.

If you want some help pondering, try this: Sketch some arbitrary function $f(\tau)$ on a set of axes at the bottom of a piece of paper—some squiggle of a curve like



will do nicely—then on a separate set of axes directly above the first, sketch the corresponding slope function $df/d\tau$. Mark two points a and b on the common horizontal axis; then on the upper, $df/d\tau$ plot, shade the integration area under the curve. Now consider (7.2) in light of your sketch.

There. Does the idea not dawn?

Another way to see the truth of the formula begins by canceling its $(1/d\tau) d\tau$ to obtain the form $\int_{\tau=a}^b df = f(\tau)|_a^b$. If this way works better for you, fine; but make sure that you understand it the other way, too.

Table 7.1: Basic derivatives for the antiderivative.

$$\int_a^b \frac{df}{d\tau} d\tau = f(\tau)|_a^b$$

$$\tau^{a-1} = \frac{d}{d\tau} \left(\frac{\tau^a}{a} \right), \quad a \neq 0$$

$$\frac{1}{\tau} = \frac{d}{d\tau} \ln \tau, \quad \ln 1 = 0$$

$$\exp \tau = \frac{d}{d\tau} \exp \tau, \quad \exp 0 = 1$$

$$\cos \tau = \frac{d}{d\tau} \sin \tau, \quad \sin 0 = 0$$

$$\sin \tau = \frac{d}{d\tau} (-\cos \tau), \quad \cos 0 = 1$$

is urged to pause now and ponder the formula thoroughly until he feels reasonably confident that indeed he does grasp it and the important idea it represents. One is unlikely to do much higher mathematics without this formula.

As an example of the formula's use, consider that because $(d/d\tau)(\tau^3/6) = \tau^2/2$, it follows that

$$\int_2^x \frac{\tau^2 d\tau}{2} = \int_2^x \frac{d}{d\tau} \left(\frac{\tau^3}{6} \right) d\tau = \frac{\tau^3}{6} \Big|_2^x = \frac{x^3 - 8}{6}.$$

Gathering elements from (4.21) and from Tables 5.2 and 5.3, Table 7.1 lists a handful of the simplest, most useful derivatives for antiderivative use. Section 9.1 speaks further of the antiderivative.

7.3 Operators, linearity and multiple integrals

This section presents the operator concept, discusses linearity and its consequences, treats the transitivity of the summational and integrodifferential operators, and introduces the multiple integral.

7.3.1 Operators

An *operator* is a mathematical agent which combines several values of a function.

Such a definition, unfortunately, is extraordinarily unilluminating to those who do not already know what it means. A better way to introduce the operator is by giving examples. Operators include $+$, $-$, multiplication, division, \sum , \prod , \int and ∂ . The essential action of an operator is to take several values of a function and combine them in some way. For example, \prod is an operator in

$$\prod_{j=1}^5 (2j-1) = (1)(3)(5)(7)(9) = 0\text{x}3\text{B}1.$$

Notice that the operator has acted to remove the variable j from the expression $2j-1$. The j appears on the equation's left side but not on its right. The operator has used the variable up. Such a variable, used up by an operator, is a *dummy variable*, as encountered earlier in § 2.3.

7.3.2 A formalism

But then how are $+$ and $-$ operators? They don't use any dummy variables up, do they? Well, it depends on how you look at it. Consider the sum $S = 3 + 5$. One can write this as

$$S = \sum_{k=0}^1 f(k),$$

where

$$f(k) \equiv \begin{cases} 3 & \text{if } k = 0, \\ 5 & \text{if } k = 1, \\ \text{undefined} & \text{otherwise.} \end{cases}$$

Then,

$$S = \sum_{k=0}^1 f(k) = f(0) + f(1) = 3 + 5 = 8.$$

By such admittedly excessive formalism, the $+$ operator can indeed be said to use a dummy variable up. The point is that $+$ is in fact an operator just like the others.

Another example of the kind:

$$\begin{aligned}
 D &= g(z) - h(z) + p(z) + q(z) \\
 &= g(z) - h(z) + p(z) - 0 + q(z) \\
 &= \Phi(0, z) - \Phi(1, z) + \Phi(2, z) - \Phi(3, z) + \Phi(4, z) \\
 &= \sum_{k=0}^4 (-1)^k \Phi(k, z),
 \end{aligned}$$

where

$$\Phi(k, z) \equiv \begin{cases} g(z) & \text{if } k = 0, \\ h(z) & \text{if } k = 1, \\ p(z) & \text{if } k = 2, \\ 0 & \text{if } k = 3, \\ q(z) & \text{if } k = 4, \\ \text{undefined} & \text{otherwise.} \end{cases}$$

Such unedifying formalism is essentially useless in applications, except as a vehicle for definition. Once you understand why $+$ and $-$ are operators just as \sum and \int are, you can forget the formalism. It doesn't help much.

7.3.3 Linearity

A function $f(z)$ is *linear* iff (if and only if) it has the properties

$$\begin{aligned}
 f(z_1 + z_2) &= f(z_1) + f(z_2), \\
 f(\alpha z) &= \alpha f(z), \\
 f(0) &= 0.
 \end{aligned}$$

The functions $f(z) = 3z$, $f(u, v) = 2u - v$ and $f(z) = 0$ are examples of linear functions. Nonlinear functions include⁶ $f(z) = z^2$, $f(u, v) = \sqrt{uv}$, $f(t) = \cos \omega t$, $f(z) = 3z + 1$ and even $f(z) = 1$.

⁶If $3z + 1$ is a *linear expression*, then how is not $f(z) = 3z + 1$ a *linear function*? Answer: it is partly a matter of purposeful definition, partly of semantics. The equation $y = 3x + 1$ plots a line, so the expression $3z + 1$ is literally “linear” in this sense; but the definition has more purpose to it than merely this. When you see the linear expression $3z + 1$, think $3z + 1 = 0$, then $g(z) = 3z = -1$. The $g(z) = 3z$ is linear; the -1 is the constant value it targets. That's the sense of it.

An operator L is linear iff it has the properties

$$\begin{aligned} L(f_1 + f_2) &= Lf_1 + Lf_2, \\ L(\alpha f) &= \alpha Lf, \\ L(0) &= 0. \end{aligned}$$

The operators \sum , \int , $+$, $-$ and ∂ are examples of linear operators. For instance,⁷

$$\frac{d}{dz}[f_1(z) + f_2(z)] = \frac{df_1}{dz} + \frac{df_2}{dz}.$$

Nonlinear operators include multiplication, division and the various trigonometric functions, among others.

7.3.4 Summational and integrodifferential transitivity

Consider the sum

$$S_1 = \sum_{k=a}^b \left[\sum_{j=p}^q \frac{x^k}{j!} \right].$$

This is a sum of the several values of the expression $x^k/j!$, evaluated at every possible pair (j, k) in the indicated domain. Now consider the sum

$$S_2 = \sum_{j=p}^q \left[\sum_{k=a}^b \frac{x^k}{j!} \right].$$

This is evidently a sum of the same values, only added in a different order. Apparently $S_1 = S_2$. Reflection along these lines must soon lead the reader to the conclusion that, in general,

$$\sum_k \sum_j f(j, k) = \sum_j \sum_k f(j, k).$$

Now consider that an integral is just a sum of many elements, and that a derivative is just a difference of two elements. Integrals and derivatives must

⁷You don't see d in the list of linear operators? But d in this context is really just another way of writing ∂ , so, yes, d is linear, too. See § 4.4.2.

then have the same transitive property discrete sums have. For example,

$$\begin{aligned}\int_{v=-\infty}^{\infty} \int_{u=a}^b f(u, v) du dv &= \int_{u=a}^b \int_{v=-\infty}^{\infty} f(u, v) dv du; \\ \int \sum_k f_k(v) dv &= \sum_k \int f_k(v) dv; \\ \frac{\partial}{\partial v} \int f du &= \int \frac{\partial f}{\partial v} du.\end{aligned}$$

In general,

$$L_v L_u f(u, v) = L_u L_v f(u, v), \quad (7.3)$$

where L is any of the linear operators \sum , \int or ∂ .

Some convergent summations, like

$$\sum_{k=0}^{\infty} \sum_{j=0}^1 \frac{(-)^j}{2k+j+1},$$

diverge once reordered, as

$$\sum_{j=0}^1 \sum_{k=0}^{\infty} \frac{(-)^j}{2k+j+1}.$$

One cannot blithely swap operators here. This is not because swapping is wrong, but rather because the inner sum after the swap diverges, hence the outer sum after the swap has no concrete summand on which to work. (*Why* does the inner sum after the swap diverge? Answer: $1 + 1/3 + 1/5 + \dots = [1] + [1/3 + 1/5] + [1/7 + 1/9 + 1/0xB + 1/0xD] + \dots > 1[1/4] + 2[1/8] + 4[1/0x10] + \dots = 1/4 + 1/4 + 1/4 + \dots$. See also § 8.10.5.) For a more twisted example of the same phenomenon, consider⁸

$$1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \dots = \left(1 - \frac{1}{2} - \frac{1}{4}\right) + \left(\frac{1}{3} - \frac{1}{6} - \frac{1}{8}\right) + \dots,$$

which associates two negative terms with each positive, but still seems to omit no term. Paradoxically, then,

$$\begin{aligned}1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \dots &= \left(\frac{1}{2} - \frac{1}{4}\right) + \left(\frac{1}{6} - \frac{1}{8}\right) + \dots \\ &= \frac{1}{2} - \frac{1}{4} + \frac{1}{6} - \frac{1}{8} + \dots \\ &= \frac{1}{2} \left(1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \dots\right),\end{aligned}$$

⁸[1, § 1.2.3]

or so it would seem, but cannot be, for it claims falsely that the sum is half itself. A better way to have handled the example might have been to write the series as

$$\lim_{n \rightarrow \infty} \left\{ 1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \cdots + \frac{1}{2n-1} - \frac{1}{2n} \right\}$$

in the first place, thus explicitly specifying equal numbers of positive and negative terms.⁹ So specifying would have prevented the error.

The *conditional convergence*¹⁰ of the last paragraph, which can occur in integrals as well as in sums, seldom poses much of a dilemma in practice. One can normally swap summational and integrodifferential operators with little worry. The reader however should at least be aware that conditional convergence troubles can arise where a summand or integrand varies in sign or phase.

7.3.5 Multiple integrals

Consider the function

$$f(u, v) = \frac{u^2}{v}.$$

Such a function would not be plotted as a curved line in a plane, but rather as a curved *surface* in a three-dimensional space. Integrating the function seeks not the area under the curve but rather the volume under the surface:

$$V = \int_{u_1}^{u_2} \int_{v_1}^{v_2} \frac{u^2}{v} dv du.$$

This is a *double integral*. Inasmuch as it can be written in the form

$$\begin{aligned} V &= \int_{u_1}^{u_2} g(u) du, \\ g(u) &\equiv \int_{v_1}^{v_2} \frac{u^2}{v} dv, \end{aligned}$$

⁹Some students of professional mathematics would assert that the false conclusion was reached through lack of rigor. Well, maybe. This writer however does not feel sure that *rigor* is quite the right word for what was lacking here. Professional mathematics does bring an elegant notation and a set of formalisms which serve ably to spotlight certain limited kinds of blunders, but these are blunders no less by the applied approach. The stalwart Leonhard Euler—arguably the greatest series-smith in mathematical history—wielded his heavy analytical hammer in thunderous strokes before professional mathematics had conceived the notation or the formalisms. If the great Euler did without, then you and I might not always be forbidden to follow his robust example.

On the other hand, the professional approach is worth study if you have the time. Recommended introductions include [24], preceded if necessary by [18] and/or [1, Ch. 1].

¹⁰[24, § 16]

its effect is to cut the area under the surface into flat, upright slices, then the slices crosswise into tall, thin towers. The towers are integrated over v to constitute the slice, then the slices over u to constitute the volume.

In light of § 7.3.4, evidently nothing prevents us from swapping the integrations: u first, then v . Hence

$$V = \int_{v_1}^{v_2} \int_{u_1}^{u_2} \frac{u^2}{v} du dv.$$

And indeed this makes sense, doesn't it? What difference does it make whether we add the towers by rows first then by columns, or by columns first then by rows? The total volume is the same in any case.

Double integrations arise very frequently in applications. Triple integrations arise about as often. For instance, if $\mu(\mathbf{r}) = \mu(x, y, z)$ represents the variable mass density of some soil,¹¹ then the total soil mass in some rectangular volume is

$$M = \int_{x_1}^{x_2} \int_{y_1}^{y_2} \int_{z_1}^{z_2} \mu(x, y, z) dz dy dx.$$

As a concise notational convenience, the last is often written

$$M = \int_V \mu(\mathbf{r}) d\mathbf{r},$$

where the V stands for “volume” and is understood to imply a triple integration. Similarly for the double integral,

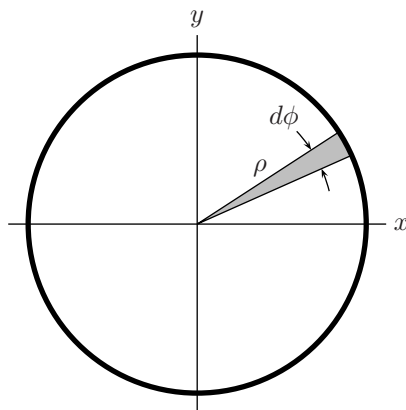
$$V = \int_S f(\boldsymbol{\rho}) d\boldsymbol{\rho},$$

where the S stands for “surface” and is understood to imply a double integration.

Even more than three nested integrations are possible. If we integrated over time as well as space, the integration would be fourfold. A spatial Fourier transform ([section not yet written]) implies a triple integration; and its inverse, another triple: a sixfold integration altogether. Manifold nesting of integrals is thus not just a theoretical mathematical topic; it arises in sophisticated real-world engineering models. The topic concerns us here for this reason.

¹¹Conventionally the Greek letter ρ not μ is used for density, but it happens that we need the letter ρ for a different purpose later in the paragraph.

Figure 7.4: The area of a circle.



7.4 Areas and volumes

By composing and solving appropriate integrals, one can calculate the perimeters, areas and volumes of interesting common shapes and solids.

7.4.1 The area of a circle

Figure 7.4 depicts an element of a circle's area. The element has wedge shape, but inasmuch as the wedge is infinitesimally narrow, the wedge is indistinguishable from a triangle of base length $\rho d\phi$ and height ρ . The area of such a triangle is $A_{\text{triangle}} = \rho^2 d\phi/2$. Integrating the many triangles, we find the circle's area to be

$$A_{\text{circle}} = \int_{\phi=-\pi}^{\pi} A_{\text{triangle}} = \int_{-\pi}^{\pi} \frac{\rho^2 d\phi}{2} = \frac{2\pi\rho^2}{2}. \quad (7.4)$$

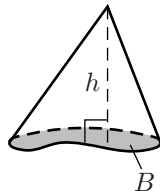
(The numerical value of 2π —the circumference or perimeter of the unit circle—we have not calculated yet. We shall calculate it in § 8.11.)

7.4.2 The volume of a cone

One can calculate the volume of any cone (or pyramid) if one knows its base area B and its altitude h measured normal¹² to the base. Refer to Fig. 7.5.

¹²Normal here means “at right angles.”

Figure 7.5: The volume of a cone.



A cross-section of a cone, cut parallel to the cone's base, has the same shape the base has but a different scale. If coordinates are chosen such that the altitude h runs in the \hat{z} direction with $z = 0$ at the cone's vertex, then the cross-sectional area is evidently¹³ $(B)(z/h)^2$. For this reason, the cone's volume is

$$V_{\text{cone}} = \int_0^h (B) \left(\frac{z}{h}\right)^2 dz = \frac{B}{h^2} \int_0^h z^2 dz = \frac{B}{h^2} \left(\frac{h^3}{3}\right) = \frac{Bh}{3}. \quad (7.5)$$

7.4.3 The surface area and volume of a sphere

Of a sphere, Fig. 7.6, one wants to calculate both the surface area and the volume. For the surface area, the sphere's surface is sliced vertically down the z axis into narrow constant- ϕ tapered strips (each strip broadest at the sphere's equator, tapering to points at the sphere's $\pm z$ poles) and horizontally across the z axis into narrow constant- θ rings, as in Fig. 7.7. A surface element so produced (seen as shaded in the latter figure) evidently has the area

$$dS = (r d\theta)(\rho d\phi) = r^2 \sin \theta d\theta d\phi.$$

¹³The fact may admittedly not be evident to the reader at first glance. If it is not yet evident to you, then ponder Fig. 7.5 a moment. Consider what it means to cut parallel to a cone's base a cross-section of the cone, and how cross-sections cut nearer a cone's vertex are smaller though the same shape. What if the base were square? Would the cross-sectional area not be $(B)(z/h)^2$ in that case? What if the base were a right triangle with equal legs—in other words, half a square? What if the base were some other strange shape like the base depicted in Fig. 7.5? Could such a strange shape not also be regarded as a definite, well characterized part of a square? (With a pair of scissors one can cut any shape from a square piece of paper, after all.) Thinking along such lines must soon lead one to the insight that the parallel-cut cross-sectional area of a cone can be nothing other than $(B)(z/h)^2$, regardless of the base's shape.

Figure 7.6: A sphere.

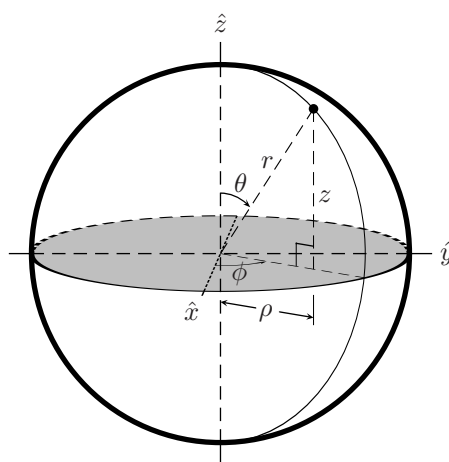
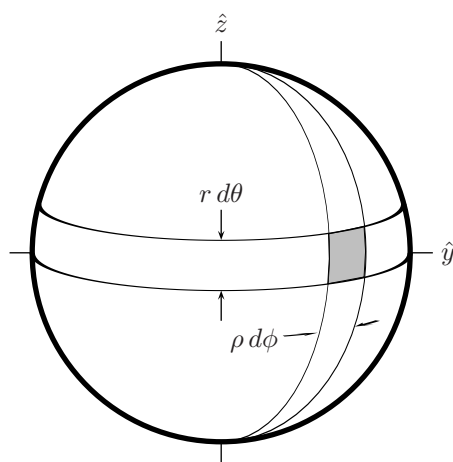


Figure 7.7: An element of the sphere's surface (see Fig. 7.6).



The sphere's total surface area then is the sum of all such elements over the sphere's entire surface:

$$\begin{aligned}
 S_{\text{sphere}} &= \int_{\phi=-\pi}^{\pi} \int_{\theta=0}^{\pi} dS \\
 &= \int_{\phi=-\pi}^{\pi} \int_{\theta=0}^{\pi} r^2 \sin \theta \, d\theta \, d\phi \\
 &= r^2 \int_{\phi=-\pi}^{\pi} [-\cos \theta]_0^{\pi} d\phi \\
 &= r^2 \int_{\phi=-\pi}^{\pi} [2] d\phi \\
 &= 4\pi r^2,
 \end{aligned} \tag{7.6}$$

where we have used the fact from Table 7.1 that $\sin \tau = (d/d\tau)(-\cos \tau)$.

Having computed the sphere's surface area, one can find its volume just as § 7.4.1 has found a circle's area—except that instead of dividing the circle into many narrow triangles, one divides the sphere into many narrow *cones*, each cone with base area dS and altitude r , with the vertices of all the cones meeting at the sphere's center. Per (7.5), the volume of one such cone is $V_{\text{cone}} = r \, dS/3$. Hence,

$$V_{\text{sphere}} = \oint_S V_{\text{cone}} = \oint_S \frac{r \, dS}{3} = \frac{r}{3} \oint_S dS = \frac{r}{3} S_{\text{sphere}},$$

where the useful symbol

$$\oint_S$$

indicates *integration over a closed surface*. In light of (7.6), the total volume is

$$V_{\text{sphere}} = \frac{4\pi r^3}{3}. \tag{7.7}$$

(One can compute the same spherical volume more prosaically, without reference to cones, by writing $dV = r^2 \sin \theta \, dr \, d\theta \, d\phi$ then integrating $\int_V dV$. The derivation given above, however, is preferred because it lends the additional insight that a sphere can sometimes be viewed as a great cone rolled up about its own vertex. The circular area derivation of § 7.4.1 lends an analogous insight: that a circle can sometimes be viewed as a great triangle rolled up about *its* own vertex.)

7.5 Checking an integration

Dividing $0x46B/0xD = 0x57$ with a pencil, how does one check the result?¹⁴ Answer: by multiplying $(0x57)(0xD) = 0x46B$. Multiplication inverts division. Easier than division, multiplication provides a quick, reliable check.

Likewise, integrating

$$\int_a^b \frac{\tau^2}{2} d\tau = \frac{b^3 - a^3}{6}$$

with a pencil, how does one check the result? Answer: by differentiating

$$\left[\frac{\partial}{\partial b} \left(\frac{b^3 - a^3}{6} \right) \right]_{b=\tau} = \frac{\tau^2}{2}.$$

Differentiation inverts integration. Easier than integration, differentiation like multiplication provides a quick, reliable check.

More formally, according to (7.2),

$$S \equiv \int_a^b \frac{df}{d\tau} d\tau = f(b) - f(a). \quad (7.8)$$

Differentiating (7.8) with respect to b and a ,

$$\begin{aligned} \left. \frac{\partial S}{\partial b} \right|_{b=\tau} &= \frac{df}{d\tau}, \\ \left. \frac{\partial S}{\partial a} \right|_{a=\tau} &= -\frac{df}{d\tau}. \end{aligned} \quad (7.9)$$

Either line of (7.9) can be used to check an integration. Evaluating (7.8) at $b = a$ yields

$$S|_{b=a} = 0, \quad (7.10)$$

which can be used to check further.¹⁵

As useful as (7.9) and (7.10) are, they nevertheless serve only integrals with variable limits. They are of little use to check *definite integrals*

¹⁴Admittedly, few readers will ever have done much such multidigit *hexadecimal* arithmetic with a pencil, but, hey, go with it. In decimal, it's $1131/13 = 87$.

Actually, hexadecimal is just proxy for binary (see Appendix A), and long division in straight binary is kind of fun. If you have never tried it, you might. It is simpler than decimal or hexadecimal division, and it's how computers divide. The insight gained is worth the trial.

¹⁵Using (7.10) to check the example, $(b^3 - a^3)/6|_{b=a} = 0$.

like (9.14) below, which lack variable limits to differentiate. However, many or most integrals one meets in practice have or can be given variable limits. Equations (7.9) and (7.10) do serve such *indefinite integrals*.

It is a rare irony of mathematics that, although numerically differentiation is indeed harder than integration, analytically precisely the opposite is true. Analytically, differentiation is the easier. So far the book has introduced only easy integrals, but Ch. 9 will bring much harder ones. Even experienced mathematicians are apt to err in analyzing these. Reversing an integration by taking an easy derivative is thus an excellent way to check a hard-earned integration result.

7.6 Contour integration

To this point we have considered only integrations in which the variable of integration advances in a straight line from one point to another: for instance, $\int_a^b f(\tau) d\tau$, in which the function $f(\tau)$ is evaluated at $\tau = a, a + d\tau, a + 2d\tau, \dots, b$. The integration variable is a real-valued scalar which can do nothing but make a straight line from a to b .

Such is not the case when the integration variable is a vector. Consider the integral

$$S = \int_{\mathbf{r}=\hat{\mathbf{x}}\rho}^{\hat{\mathbf{y}}\rho} (x^2 + y^2) d\ell,$$

where $d\ell$ is the infinitesimal length of a step along the path of integration. What does this integral mean? Does it mean to integrate from $\mathbf{r} = \hat{\mathbf{x}}\rho$ to $\mathbf{r} = 0$, then from there to $\mathbf{r} = \hat{\mathbf{y}}\rho$? Or does it mean to integrate along the arc of Fig. 7.8? The two paths of integration begin and end at the same points, but they differ in between, and the integral certainly does not come out the same both ways. Yet many other paths of integration from $\hat{\mathbf{x}}\rho$ to $\hat{\mathbf{y}}\rho$ are possible, not just these two.

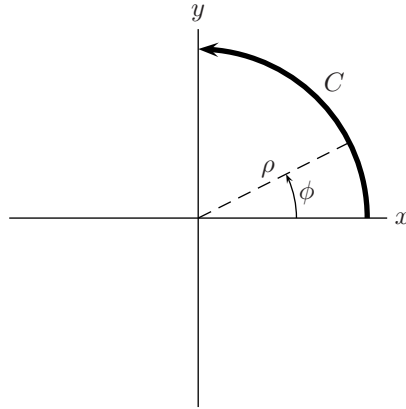
Because multiple paths are possible, we must be more specific:

$$S = \int_C (x^2 + y^2) d\ell,$$

where C stands for “contour” and means in this example the specific contour of Fig. 7.8. In the example, $x^2 + y^2 = \rho^2$ (by the Pythagorean theorem) and $d\ell = \rho d\phi$, so

$$S = \int_C \rho^2 d\ell = \int_0^{2\pi/4} \rho^3 d\phi = \frac{2\pi}{4} \rho^3.$$

Figure 7.8: A contour of integration.



In the example the contour is open, but closed contours which begin and end at the same point are also possible, indeed common. The useful symbol

$$\oint$$

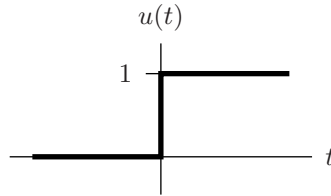
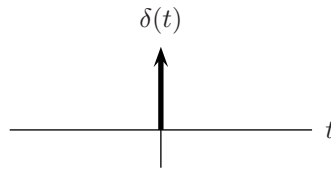
indicates *integration over a closed contour*. It means that the contour ends where it began: the loop is closed. The contour of Fig. 7.8 would be closed, for instance, if it continued to $\mathbf{r} = 0$ and then back to $\mathbf{r} = \hat{\mathbf{x}}\rho$.

Besides applying where the variable of integration is a vector, contour integration applies equally where the variable of integration is a complex scalar. In the latter case some interesting mathematics emerge, as we shall see in §§ 8.8 and 9.5.

7.7 Discontinuities

The polynomials and trigonometrics studied to this point in the book offer flexible means to model many physical phenomena of interest, but one thing they do not model gracefully is the simple discontinuity. Consider a mechanical valve opened at time $t = t_o$. The flow $x(t)$ past the valve is

$$x(t) = \begin{cases} 0, & t < t_o; \\ x_o, & t > t_o. \end{cases}$$

Figure 7.9: The Heaviside unit step $u(t)$.Figure 7.10: The Dirac delta $\delta(t)$.

One can write this more concisely in the form

$$x(t) = u(t - t_o)x_o,$$

where $u(t)$ is the *Heaviside unit step*,

$$u(t) \equiv \begin{cases} 0, & t < 0; \\ 1, & t > 0; \end{cases} \quad (7.11)$$

plotted in Fig. 7.9.

The derivative of the Heaviside unit step is the curious *Dirac delta*

$$\delta(t) \equiv \frac{d}{dt}u(t), \quad (7.12)$$

plotted in Fig. 7.10. This function is zero everywhere except at $t = 0$, where it is infinite, with the property that

$$\int_{-\infty}^{\infty} \delta(t) dt = 1, \quad (7.13)$$

and the interesting consequence that

$$\int_{-\infty}^{\infty} \delta(t - t_o) f(t) dt = f(t_o) \quad (7.14)$$

for any function $f(t)$. (Equation 7.14 is the *sifting property* of the Dirac delta.)¹⁶

The Dirac delta is defined for vectors, too, such that

$$\int_V \delta(\mathbf{r}) d\mathbf{r} = 1. \quad (7.15)$$

¹⁶It seems inadvisable for the narrative to digress at this point to explore $u(z)$ and $\delta(z)$, the unit step and delta of a complex argument, although by means of Fourier analysis ([chapter not yet written]) or by conceiving the Dirac delta as an infinitely narrow Gaussian pulse ([chapter not yet written]) it could perhaps do so. The book has more pressing topics to treat. For the book's present purpose the interesting action of the two functions is with respect to the real argument t .

In the author's country at least, a sort of debate seems to have run for decades between professional and applied mathematicians over the Dirac delta $\delta(t)$. Some professional mathematicians seem to have objected that $\delta(t)$ is not a function, inasmuch as it lacks certain properties common to functions as they define them [32, § 2.4][10]. From the applied point of view the objection is admittedly a little hard to understand, until one realizes that it is more a dispute over methods and definitions than over facts. What the professionals seem to be saying is that $\delta(t)$ does not fit as neatly as they would like into the abstract mathematical framework they had established for functions in general before Paul Dirac came along in 1930 [46, "Paul Dirac," 05:48, 25 May 2006] and slapped his disruptive $\delta(t)$ down on the table. The objection is not so much that $\delta(t)$ is not allowed as it is that professional mathematics for years after 1930 lacked a fully coherent theory for it.

It's a little like the six-fingered man in Goldman's *The Princess Bride* [17]. If I had established a definition of "nobleman" which subsumed "human," whose relevant traits in my definition included five fingers on each hand, when the six-fingered Count Rugen appeared on the scene, then you would expect me to adapt my definition, wouldn't you? By my preëxisting definition, strictly speaking, the six-fingered count is "not a nobleman"; but such exclusion really tells one more about flaws in the definition than it does about the count.

Whether the professional mathematician's definition of the *function* is flawed, of course, is not for this writer to judge. Even if not, however, the fact of the Dirac delta dispute, coupled with the difficulty we applied mathematicians experience in trying to understand the reason the dispute even exists, has unfortunately surrounded the Dirac delta with a kind of mysterious aura, an elusive sense that $\delta(t)$ hides subtle mysteries—when what it really hides is an internal discussion of words and means among the professionals. The professionals who had established the theoretical framework before 1930 justifiably felt reluctant to throw the whole framework away because some scientists and engineers like us came along one day with a useful new function which didn't quite fit, but that was the professionals' problem not ours. To us the Dirac delta $\delta(t)$ is just a function. The internal discussion of words and means, we leave to the professionals, who know whereof they speak.

7.8 Remarks (and exercises)

The concept of the integral is relatively simple once grasped, but its implications are broad, deep and hard. This chapter is short. One reason introductory calculus texts run so long is that they include many, many pages of integration examples and exercises. The reader who desires a gentler introduction to the integral might consult among others the textbook the chapter's introduction has recommended.

Even if this book is not an instructional textbook, it seems not meet that it should include no exercises at all here. Here are a few. Some of them do need material from later chapters, so you should not expect to be able to complete them all now. The harder ones are marked with *asterisks. Work the exercises if you like.

1. Evaluate (a) $\int_0^x \tau d\tau$; (b) $\int_0^x \tau^2 d\tau$. (Answer: $x^2/2$; $x^3/3$.)
2. Evaluate (a) $\int_1^x (1/\tau^2) d\tau$; (b) $\int_a^x 3\tau^{-2} d\tau$; (c) $\int_a^x C\tau^n d\tau$; (d) $\int_0^x (a_2\tau^2 + a_1\tau) d\tau$; *(e) $\int_1^x (1/\tau) d\tau$.
3. *Evaluate (a) $\int_0^x \sum_{k=0}^{\infty} \tau^k d\tau$; (b) $\sum_{k=0}^{\infty} \int_0^x \tau^k d\tau$; (c) $\int_0^x \sum_{k=0}^{\infty} (\tau^k/k!) d\tau$.
4. Evaluate $\int_0^x \exp \alpha \tau d\tau$.
5. Evaluate (a) $\int_{-2}^5 (3\tau^2 - 2\tau^3) d\tau$; (b) $\int_5^{-2} (3\tau^2 - 2\tau^3) d\tau$. Work the exercise by hand in hexadecimal and give the answer in hexadecimal.
6. Evaluate $\int_1^{\infty} (3/\tau^2) d\tau$.
7. *Evaluate the integral of the example of § 7.6 along the alternate contour suggested there, from $\hat{\mathbf{x}}\rho$ to 0 to $\hat{\mathbf{y}}\rho$.
8. Evaluate (a) $\int_0^x \cos \omega \tau d\tau$; (b) $\int_0^x \sin \omega \tau d\tau$; *(c)¹⁷ $\int_0^x \tau \sin \omega \tau d\tau$.
9. *Evaluate¹⁸ (a) $\int_1^x \sqrt{1+2\tau} d\tau$; (b) $\int_x^a [(\cos \sqrt{\tau})/\sqrt{\tau}] d\tau$.
10. *Evaluate¹⁹ (a) $\int_0^x [1/(1+\tau^2)] d\tau$ (answer: $\arctan x$); (b) $\int_0^x [(4+i3)/\sqrt{2-3\tau^2}] d\tau$ (hint: the answer involves another inverse trigonometric).
11. **Evaluate (a) $\int_{-\infty}^x \exp[-\tau^2/2] d\tau$; (b) $\int_{-\infty}^{\infty} \exp[-\tau^2/2] d\tau$.

¹⁷[37, § 8-2]

¹⁸[37, § 5-6]

¹⁹[37, back endpaper]

The last exercise in particular requires some experience to answer. Moreover, it requires a developed sense of applied mathematical style to put the answer in a pleasing form (the right form for part b is very different from that for part a). Some of the easier exercises, of course, you should be able to work right now.

The point of the exercises is to illustrate how hard integrals can be to solve, and in fact how easy it is to come up with an integral which no one really knows how to solve very well. Some solutions to the same integral are better than others (easier to manipulate, faster to numerically calculate, etc.) yet not even the masters can solve them all in practical ways. On the other hand, integrals which arise in practice often can be solved very well with sufficient cleverness—and the more cleverness you develop, the more such integrals you can solve. The ways to solve them are myriad. The mathematical art of solving diverse integrals is well worth cultivating.

Chapter 9 introduces some of the basic, most broadly useful integral-solving techniques. Before addressing techniques of integration, however, as promised earlier we turn our attention in Chapter 8 back to the derivative, applied in the form of the Taylor series.

Chapter 8

The Taylor series

The Taylor series is a power series which fits a function in a limited domain neighborhood. Fitting a function in such a way brings two advantages:

- it lets us take derivatives and integrals in the same straightforward way (4.20) we take them with any power series; and
- it implies a simple procedure to calculate the function numerically.

This chapter introduces the Taylor series and some of its incidents. It also derives Cauchy's integral formula. The chapter's early sections prepare the ground for the treatment of the Taylor series proper in § 8.3.¹

¹Because even at the applied level the proper derivation of the Taylor series involves mathematical induction, analytic continuation and the matter of convergence domains, no balance of rigor the chapter might strike seems wholly satisfactory. The chapter errs maybe toward too much rigor; for, with a little less, most of §§ 8.1, 8.2, 8.4 and 8.6 would cease to be necessary. For the impatient, to read only the following sections might not be an unreasonable way to shorten the chapter: §§ 8.3, 8.5, 8.8, 8.9 and 8.11, plus the introduction of § 8.1.

From another point of view, the chapter errs maybe toward too little rigor. Some pretty constructs of pure mathematics serve the Taylor series and Cauchy's integral formula. However, such constructs drive the applied mathematician on too long a detour. The chapter as written represents the most nearly satisfactory compromise the writer has been able to strike.

8.1 The power-series expansion of $1/(1-z)^{n+1}$

Before approaching the Taylor series proper in § 8.3, we shall find it both interesting and useful to demonstrate that

$$\frac{1}{(1-z)^{n+1}} = \sum_{k=0}^{\infty} \binom{n+k}{n} z^k \quad (8.1)$$

for $n \geq 0$, $|z| < 1$. The demonstration comes in three stages. Of the three, it is the second stage (§ 8.1.2) which actually proves (8.1). The first stage (§ 8.1.1) comes up with the formula for the second stage to prove. The third stage (§ 8.1.3) establishes the sum's convergence.

8.1.1 The formula

In § 2.6.4 we found that

$$\frac{1}{1-z} = \sum_{k=0}^{\infty} z^k = 1 + z + z^2 + z^3 + \dots$$

for $|z| < 1$. What about $1/(1-z)^2$, $1/(1-z)^3$, $1/(1-z)^4$, and so on? By the long-division procedure of Table 2.4, one can calculate the first few terms of $1/(1-z)^2$ to be

$$\frac{1}{(1-z)^2} = \frac{1}{1-2z+z^2} = 1 + 2z + 3z^2 + 4z^3 + \dots$$

whose coefficients 1, 2, 3, 4, ... happen to be the numbers down the first diagonal of Pascal's triangle (Fig. 4.2 on page 73; see also Fig. 4.1). Dividing $1/(1-z)^3$ seems to produce the coefficients 1, 3, 6, 10, ... down the second diagonal; dividing $1/(1-z)^4$, the coefficients down the third. A curious pattern seems to emerge, worth investigating more closely. The pattern recommends the conjecture (8.1).

To motivate the conjecture a bit more formally (though without actually proving it yet), suppose that $1/(1-z)^{n+1}$, $n \geq 0$, is expandable in the power series

$$\frac{1}{(1-z)^{n+1}} = \sum_{k=0}^{\infty} a_{nk} z^k, \quad (8.2)$$

where the a_{nk} are coefficients to be determined. Multiplying by $1-z$, we have that

$$\frac{1}{(1-z)^n} = \sum_{k=0}^{\infty} (a_{nk} - a_{n(k-1)}) z^k.$$

This is to say that

$$a_{(n-1)k} = a_{nk} - a_{n(k-1)},$$

or in other words that

$$a_{n(k-1)} + a_{(n-1)k} = a_{nk}. \quad (8.3)$$

Thinking of Pascal's triangle, (8.3) reminds one of (4.5), transcribed here in the symbols

$$\binom{m-1}{j-1} + \binom{m-1}{j} = \binom{m}{j}, \quad (8.4)$$

except that (8.3) is not $a_{(m-1)(j-1)} + a_{(m-1)j} = a_{mj}$.

Various changes of variable are possible to make (8.4) better match (8.3). We might try at first a few false ones, but eventually the change

$$\begin{aligned} n+k &\leftarrow m, \\ k &\leftarrow j, \end{aligned}$$

recommends itself. Thus changing in (8.4) gives

$$\binom{n+k-1}{k-1} + \binom{n+k-1}{k} = \binom{n+k}{k}.$$

Transforming according to the rule (4.3), this is

$$\binom{n+(k-1)}{n} + \binom{(n-1)+k}{n-1} = \binom{n+k}{n}, \quad (8.5)$$

which fits (8.3) perfectly. Hence we conjecture that

$$a_{nk} = \binom{n+k}{n}, \quad (8.6)$$

which coefficients, applied to (8.2), yield (8.1).

Equation (8.1) is thus suggestive. It works at least for the important case of $n = 0$; this much is easy to test. In light of (8.3), it seems to imply a relationship between the $1/(1 - z)^{n+1}$ series and the $1/(1 - z)^n$ series for any n . But *to seem* is not *to be*. At this point, all we can say is that (8.1) seems right. We shall establish that it is right in the next subsection.

8.1.2 The proof by induction

Equation (8.1) is proved by induction as follows. Consider the sum

$$S_n \equiv \sum_{k=0}^{\infty} \binom{n+k}{n} z^k. \quad (8.7)$$

Multiplying by $1 - z$ yields

$$(1 - z)S_n = \sum_{k=0}^{\infty} \left[\binom{n+k}{n} - \binom{n+(k-1)}{n} \right] z^k.$$

Per (8.5), this is

$$(1 - z)S_n = \sum_{k=0}^{\infty} \binom{(n-1)+k}{n-1} z^k. \quad (8.8)$$

Now suppose that (8.1) is true for $n = i - 1$ (where i denotes an integer rather than the imaginary unit):

$$\frac{1}{(1 - z)^i} = \sum_{k=0}^{\infty} \binom{(i-1)+k}{i-1} z^k. \quad (8.9)$$

In light of (8.8), this means that

$$\frac{1}{(1 - z)^i} = (1 - z)S_i.$$

Dividing by $1 - z$,

$$\frac{1}{(1 - z)^{i+1}} = S_i.$$

Applying (8.7),

$$\frac{1}{(1 - z)^{i+1}} = \sum_{k=0}^{\infty} \binom{i+k}{i} z^k. \quad (8.10)$$

Evidently (8.9) implies (8.10). In other words, if (8.1) is true for $n = i - 1$, then it is also true for $n = i$. Thus *by induction*, if it is true for any one n , then it is also true for all greater n .

The “if” in the last sentence is important. Like all inductions, this one needs at least one *start case* to be valid (many inductions actually need a consecutive pair of start cases). The $n = 0$ supplies the start case

$$\frac{1}{(1 - z)^{0+1}} = \sum_{k=0}^{\infty} \binom{k}{0} z^k = \sum_{k=0}^{\infty} z^k,$$

which per (2.34) we know to be true.

8.1.3 Convergence

The question remains as to the domain over which the sum (8.1) converges.² To answer the question, consider that per (4.9),

$$\binom{m}{j} = \frac{m}{m-j} \binom{m-1}{j}$$

for any integers $m > 0$ and j . With the substitution $n+k \leftarrow m$, $n \leftarrow j$, this means that

$$\binom{n+k}{n} = \frac{n+k}{k} \binom{n+(k-1)}{n},$$

or more tersely,

$$a_{nk} = \frac{n+k}{k} a_{n(k-1)},$$

where

$$a_{nk} \equiv \binom{n+k}{n}$$

are the coefficients of the power series (8.1). Rearranging factors,

$$\frac{a_{nk}}{a_{n(k-1)}} = \frac{n+k}{k} = 1 + \frac{n}{k}. \quad (8.11)$$

²The meaning of the verb *to converge* may seem clear enough from the context and from earlier references, but if explanation here helps: a series converges if and only if it approaches a specific, finite value after many terms. A more rigorous way of saying the same thing is as follows: the series

$$S = \sum_{k=0}^{\infty} \tau_k$$

converges iff (if and only if), for all possible positive constants ϵ , there exists a finite $K \geq -1$ such that

$$\left| \sum_{k=K+1}^n \tau_k \right| < \epsilon,$$

for all $n \geq K$ (of course it is also required that the τ_k be finite, but you knew that already).

The professional mathematical literature calls such convergence “uniform convergence,” distinguishing it through a test devised by Weierstrass from the weaker “pointwise convergence” [1, § 1.5.3]. The applied mathematician can profit substantially by learning the professional view in the matter, but the effect of trying to teach the professional view in a book like this would not be pleasing. Here, we avoid error by keeping a clear view of the physical phenomena the mathematics is meant to model.

It is interesting nevertheless to consider an example of an integral for which convergence is not so simple, such as Frullani’s integral of § 9.7.

Multiplying (8.11) by z^k/z^{k-1} gives the ratio

$$\frac{a_{nk}z^k}{a_{n(k-1)}z^{k-1}} = \left(1 + \frac{n}{k}\right)z,$$

which is to say that the k th term of (8.1) is $(1 + n/k)z$ times the $(k-1)$ th term. So long as the criterion³

$$\left|\left(1 + \frac{n}{k}\right)z\right| \leq 1 - \delta$$

is satisfied for all sufficiently large $k > K$ —where $0 < \delta \ll 1$ is a small positive constant—then the series evidently converges (see § 2.6.4 and eqn. 3.20). But we can bind $1 + n/k$ as close to unity as desired by making K sufficiently large, so to meet the criterion it suffices that

$$|z| < 1. \tag{8.12}$$

The bound (8.12) thus establishes a sure convergence domain for (8.1).

8.1.4 General remarks on mathematical induction

We have proven (8.1) by means of a mathematical induction. The virtue of induction as practiced in § 8.1.2 is that it makes a logically clean, airtight case for a formula. Its vice is that it conceals the subjective process which has led the mathematician to consider the formula in the first place. Once you obtain a formula somehow, maybe you can prove it by induction; but the induction probably does not help you to obtain the formula! A good inductive proof usually begins by motivating the formula proven, as in § 8.1.1.

Richard W. Hamming once said of mathematical induction,

The theoretical difficulty the student has with mathematical induction arises from the reluctance to ask seriously, “How could I prove a formula for an infinite number of cases when I know that testing a finite number of cases is not enough?” Once you

³Although one need not ask the question to understand the proof, the reader may nevertheless wonder why the simpler $|(1 + n/k)z| < 1$ is not given as a criterion. The surprising answer is that not all series $\sum \tau_k$ with $|\tau_k/\tau_{k-1}| < 1$ converge! For example, the extremely simple $\sum 1/k$ does not converge. As we see however, all series $\sum \tau_k$ with $|\tau_k/\tau_{k-1}| < 1 - \delta$ do converge. The distinction is subtle but rather important.

The really curious reader may now ask why $\sum 1/k$ does not converge. Answer: it majorizes $\int_1^x (1/\tau) d\tau = \ln x$. See (5.7) and § 8.10.

really face this question, you will understand the ideas behind mathematical induction. It is only when you grasp the problem clearly that the method becomes clear. [18, § 2.3]

Hamming also wrote,

The function of rigor is mainly critical and is seldom constructive. Rigor is the hygiene of mathematics, which is needed to protect us against careless thinking. [18, § 1.6]

The applied mathematician may tend to avoid rigor for which he finds no immediate use, but he does not disdain mathematical rigor on principle. The style lies in exercising rigor at the right level for the problem at hand. Hamming, a professional mathematician who sympathized with the applied mathematician's needs, wrote further,

Ideally, when teaching a topic the degree of rigor should follow the student's perceived need for it. . . . It is necessary to require a gradually rising level of rigor so that when faced with a real need for it you are not left helpless. As a result, [one cannot teach] a uniform level of rigor, but rather a gradually rising level. Logically, this is indefensible, but psychologically there is little else that can be done. [18, § 1.6]

Applied mathematics holds that the practice *is* defensible, on the ground that the math serves the model; but Hamming nevertheless makes a pertinent point.

Mathematical induction is a broadly applicable technique for constructing mathematical proofs. We shall not always write inductions out as explicitly in this book as we have done in the present section—often we shall leave the induction as an implicit exercise for the interested reader—but this section's example at least lays out the general pattern of the technique.

8.2 Shifting a power series' expansion point

One more question we should treat before approaching the Taylor series proper in § 8.3 concerns the shifting of a power series' expansion point. How can the expansion point of the power series

$$f(z) = \sum_{k=K}^{\infty} (a_k)(z - z_o)^k, \quad K \leq 0, \quad (8.13)$$

which may have terms of negative order, be shifted from $z = z_o$ to $z = z_1$?

The first step in answering the question is straightforward: one rewrites (8.13) in the form

$$f(z) = \sum_{k=K}^{\infty} (a_k) ([z - z_1] - [z_o - z_1])^k,$$

then changes the variables

$$\begin{aligned} w &\leftarrow \frac{z - z_1}{z_o - z_1}, \\ c_k &\leftarrow [-(z_o - z_1)]^k a_k, \end{aligned} \tag{8.14}$$

to obtain

$$f(z) = \sum_{k=K}^{\infty} (c_k) (1 - w)^k. \tag{8.15}$$

Splitting the $k < 0$ terms from the $k \geq 0$ terms in (8.15), we have that

$$\begin{aligned} f(z) &= f_-(z) + f_+(z), \\ f_-(z) &\equiv \sum_{k=0}^{-(K+1)} \frac{c_{[-(k+1)]}}{(1-w)^{k+1}}, \\ f_+(z) &\equiv \sum_{k=0}^{\infty} (c_k) (1-w)^k. \end{aligned} \tag{8.16}$$

Of the two subseries, the $f_-(z)$ is expanded term by term using (8.1), after which combining like powers of w yields the form

$$\begin{aligned} f_-(z) &= \sum_{k=0}^{\infty} q_k w^k, \\ q_k &\equiv \sum_{n=0}^{-(K+1)} (c_{[-(n+1)]}) \binom{n+k}{n}. \end{aligned} \tag{8.17}$$

The $f_+(z)$ is even simpler to expand: one need only multiply the series out term by term per (4.12), combining like powers of w to reach the form

$$\begin{aligned} f_+(z) &= \sum_{k=0}^{\infty} p_k w^k, \\ p_k &\equiv \sum_{n=k}^{\infty} (c_n) \binom{n}{k}. \end{aligned} \tag{8.18}$$

Equations (8.13) through (8.18) serve to shift a power series' expansion point, calculating the coefficients of a power series for $f(z)$ about $z = z_1$, given those of a power series about $z = z_o$. Notice that—unlike the original, $z = z_o$ power series—the new, $z = z_1$ power series has terms $(z - z_1)^k$ only for $k \geq 0$; it has no terms of negative order. At the price per (8.12) of restricting the convergence domain to $|w| < 1$, shifting the expansion point away from the pole at $z = z_o$ has resolved the $k < 0$ terms.

The method fails if $z = z_1$ happens to be a pole or other nonanalytic point of $f(z)$. The convergence domain vanishes as z_1 approaches such a forbidden point. (Examples of such forbidden points include $z = 0$ in $h[z] = 1/z$ and in $g[z] = \sqrt{z}$. See §§ 8.4 through 8.8.) Furthermore, even if z_1 does represent a fully analytic point of $f(z)$, it also must lie within the convergence domain of the original, $z = z_o$ series for the shift to be trustworthy as derived.

The attentive reader might observe that we have formally established the convergence neither of $f_-(z)$ in (8.17) nor of $f_+(z)$ in (8.18). Regarding the former convergence, that of $f_-(z)$, we have strategically framed the problem so that one needn't worry about it, running the sum in (8.13) from the finite $k = K \leq 0$ rather than from the infinite $k = -\infty$; and since according to (8.12) each term of the original $f_-(z)$ of (8.16) converges for $|w| < 1$, the reconstituted $f_-(z)$ of (8.17) safely converges in the same domain. The latter convergence, that of $f_+(z)$, is harder to establish in the abstract because that subseries has an infinite number of terms. As we shall see by pursuing a different line of argument in § 8.3, however, the $f_+(z)$ of (8.18) can be nothing other than the Taylor series about $z = z_1$ of the function $f_+(z)$ in any event, enjoying the same convergence domain any such Taylor series enjoys.⁴

8.3 Expanding functions in Taylor series

Having prepared the ground, we now stand in position to treat the Taylor series proper. The treatment begins with a question: if you had to express

⁴A rigorous argument can be constructed without appeal to § 8.3 if desired, from the ratio $n/(n - k)$ of (4.9) and its brethren, which ratio approaches unity with increasing n . A more elegant rigorous argument can be made indirectly by way of a complex contour integral. In applied mathematics, however, one does not normally try to shift the expansion point of an *unspecified* function $f(z)$, anyway. Rather, one shifts the expansion point of some concrete function like $\sin z$ or $\ln(1 - z)$. The imagined difficulty (if any) vanishes in the concrete case. Appealing to § 8.3, the important point is the one made in the narrative: $f_+(z)$ can be nothing other than the Taylor series in any event.

some function $f(z)$ by a power series

$$f(z) = \sum_{k=0}^{\infty} (a_k)(z - z_o)^k,$$

with terms of nonnegative order $k \geq 0$ only, how would you do it? The procedure of § 8.1 worked well enough in the case of $f(z) = 1/(1-z)^{n+1}$, but it is not immediately obvious that the same procedure works more generally. What if $f(z) = \sin z$, for example?⁵

Fortunately a different way to attack the power-series expansion problem is known. It works by asking the question: what power series, having terms of nonnegative order only, most resembles $f(z)$ in the immediate neighborhood of $z = z_o$? To resemble $f(z)$, the desired power series should have $a_0 = f(z_o)$; otherwise it would not have the right value at $z = z_o$. Then it should have $a_1 = f'(z_o)$ for the right slope. Then, $a_2 = f''(z_o)/2$ for the right second derivative, and so on. With this procedure,

$$f(z) = \sum_{k=0}^{\infty} \left(\left. \frac{d^k f}{dz^k} \right|_{z=z_o} \right) \frac{(z - z_o)^k}{k!}. \quad (8.19)$$

Equation (8.19) is the *Taylor series*. Where it converges, it has all the same derivatives $f(z)$ has, so if $f(z)$ is infinitely differentiable then the Taylor series is an exact representation of the function.⁶

⁵The actual Taylor series for $\sin z$ is given in § 8.9.

⁶Further proof details may be too tiresome to inflict on applied mathematical readers. However, for readers who want a little more rigor nevertheless, the argument goes briefly as follows. Consider an infinitely differentiable function $F(z)$ and its Taylor series $f(z)$ about z_o . Let $\Delta F(z) \equiv F(z) - f(z)$ be the part of $F(z)$ not representable as a Taylor series about z_o .

If $\Delta F(z)$ is the part of $F(z)$ not representable as a Taylor series, then $\Delta F(z_o)$ and all its derivatives at z_o must be identically zero (otherwise by the Taylor series formula of eqn. 8.19, one could construct a nonzero Taylor series for $\Delta F[z_o]$ from the nonzero derivatives). However, if $F(z)$ is infinitely differentiable and if *all* the derivatives of $\Delta F(z)$ are zero at $z = z_o$, then by the unbalanced definition of the derivative from § 4.4, all the derivatives must also be zero at $z = z_o \pm \epsilon$, hence also at $z = z_o \pm 2\epsilon$, and so on. This means that $\Delta F(z) = 0$. In other words, there is no part of $F(z)$ not representable as a Taylor series.

The interested reader can fill the details in, but basically that is how the more rigorous proof goes. (A more elegant rigorous proof, preferred by the professional mathematicians [25][12] but needing significant theoretical preparation, involves integrating over a complex contour about the expansion point. We will not detail that proof here.) The reason the rigorous proof is confined to a footnote is not a deprecation of rigor as such. It is a deprecation of rigor which serves little purpose in applications. Applied mathematicians

The Taylor series is not guaranteed to converge outside some neighborhood near $z = z_o$, but where it does converge it is precise.

When $z_o = 0$, the series is also called the *Maclaurin series*. By either name, the series is a construct of great importance and tremendous practical value, as we shall soon see.

8.4 Analytic continuation

As earlier mentioned in § 2.12.3, an *analytic function* is a function which is infinitely differentiable in the domain neighborhood of interest—or, maybe more appropriately for our applied purpose, a function expressible as a Taylor series in that neighborhood. As we have seen, only one Taylor series about z_o is possible for a given function $f(z)$:

$$f(z) = \sum_{k=0}^{\infty} (a_k)(z - z_o)^k.$$

However, nothing prevents one from transposing the series to a different expansion point $z = z_1$ by the method of § 8.2, except that the transposed series may there enjoy a different convergence domain. As it happens, this section's purpose finds it convenient to swap symbols $z_o \leftrightarrow z_1$, transposing rather from expansion about $z = z_1$ to expansion about $z = z_o$. In the swapped notation, so long as the expansion point $z = z_o$ lies fully within (neither outside nor right on the edge of) the $z = z_1$ series' convergence domain, the two series evidently describe the selfsame underlying analytic function.

Since an analytic function $f(z)$ is infinitely differentiable and enjoys a unique Taylor expansion $f_o(z - z_o) = f(z)$ about each point z_o in its domain, it follows that if two Taylor series $f_1(z - z_1)$ and $f_2(z - z_2)$ find even a small neighborhood $|z - z_o| < \epsilon$ which lies in the domain of both, then the two can

normally regard mathematical functions to be imprecise analogs of physical quantities of interest. Since the functions are imprecise analogs in any case, the applied mathematician is logically free implicitly *to define* the functions he uses as Taylor series in the first place; that is, to restrict the set of infinitely differentiable functions used in the model to the subset of such functions representable as Taylor series. With such an implicit definition, whether there actually exist any infinitely differentiable functions not representable as Taylor series is more or less beside the point.

In applied mathematics, the definitions serve the model, not the other way around.

(It is entertaining to consider [45, "Extremum"] the Taylor series of the function $\sin[1/x]$ —although in practice this particular function is readily expanded after the obvious change of variable $u \leftarrow 1/x$.)

both be transposed to the common $z = z_o$ expansion point. If the two are found to have the same Taylor series there, then f_1 and f_2 both represent the same function. Moreover, if a series f_3 is found whose domain overlaps that of f_2 , then a series f_4 whose domain overlaps that of f_3 , and so on, and if each pair in the chain matches at least in a small neighborhood in its region of overlap, then the whole chain of overlapping series necessarily represents the same underlying analytic function $f(z)$. The series f_1 and the series f_n represent the same analytic function even if their domains do not directly overlap at all.

This is a manifestation of the principle of *analytic continuation*. The principle holds that if two analytic functions are the same within some domain neighborhood $|z - z_o| < \epsilon$, then they are the same everywhere.⁷ Observe however that the principle fails at poles and other nonanalytic points, because the function is not differentiable there.

The result of § 8.2, which shows general power series to be expressible as Taylor series except at their poles and other nonanalytic points, extends the analytic continuation principle to cover power series in general, including power series with terms of negative order.

Now, observe: though all convergent power series are indeed analytic, one need not actually expand every analytic function in a power series. Sums, products and ratios of analytic functions are no less differentiable than the functions themselves—as also, by the derivative chain rule, is an analytic function of analytic functions. For example, where $g(z)$ and $h(z)$ are analytic, there also is $f(z) \equiv g(z)/h(z)$ analytic (except perhaps at isolated points where $h[z] = 0$). Besides, given Taylor series for $g(z)$ and $h(z)$ one can make a power series for $f(z)$ by long division if desired, so that is all right. Section 8.15 speaks further on the point.

The subject of analyticity is rightly a matter of deep concern to the professional mathematician. It is also a long wedge which drives pure and applied mathematics apart. When the professional mathematician speaks generally of a “function,” he means *any function at all*. One can construct

⁷The writer hesitates to mention that he is given to understand [39] that the domain neighborhood can technically be reduced to a domain contour of nonzero length but zero width. Having never met a significant application of this extension of the principle, the writer has neither researched the extension’s proof nor asserted its truth. He does not especially recommend that the reader worry over the point. The domain neighborhood $|z - z_o| < \epsilon$ suffices.

some pretty unreasonable functions if one wants to, such as

$$\begin{aligned} f([2k+1]2^m) &\equiv (-)^m, \\ f(z) &\equiv 0 \text{ otherwise,} \end{aligned}$$

where k and m are integers. However, neither functions like this $f(z)$ nor more subtly unreasonable functions normally arise in the modeling of physical phenomena. When such functions do arise, one transforms, approximates, reduces, replaces and/or avoids them. The full theory which classifies and encompasses—or explicitly excludes—such functions is thus of limited interest to the applied mathematician, and this book does not cover it.⁸

This does not mean that the scientist or engineer never encounters non-analytic functions. On the contrary, he encounters several, but they are not subtle: $|z|$; $\arg z$; z^* ; $\Re(z)$; $\Im(z)$; $u(t)$; $\delta(t)$. Refer to §§ 2.12 and 7.7. Such functions are nonanalytic either because they lack proper derivatives in the Argand plane according to (4.19) or because one has defined them only over a real domain.

8.5 Branch points

The function $g(z) = \sqrt{z}$ is an interesting, troublesome function. Its derivative is $dg/dz = 1/2\sqrt{z}$, so even though the function is finite at $z = 0$, its derivative is not finite there. Evidently $g(z)$ has a nonanalytic point at $z = 0$, yet the point is not a pole. What is it?

We call it a *branch point*. The defining characteristic of the branch point is that, given a function $f(z)$ with such a point at $z = z_0$, if one encircles⁹ the point once alone (that is, without also encircling some other branch point) by a closed contour in the Argand domain plane, while simultaneously tracking $f(z)$ in the Argand range plane—and if one demands that z and $f(z)$ move

⁸Many books do cover it in varying degrees, including [12][39][21] and numerous others. The foundations of the pure theory of a complex variable, though abstract, are beautiful; and though they do not comfortably fit a book like this, even an applied mathematician can profit substantially by studying them. Maybe a future release of the book will trace the pure theory's main thread in an appendix. However that may be, the pure theory is probably best appreciated after one already understands its chief conclusions. Though not for the explicit purpose of serving the pure theory, the present chapter does develop just such an understanding.

⁹For readers whose native language is not English, “to encircle” means “to surround” or “to enclose.” The verb does not require the boundary to have the shape of an actual, geometrical circle; any closed shape suffices. However, the circle is a typical shape, probably the most fitting shape to imagine when thinking of the concept abstractly.

smoothly, that neither suddenly skip from one spot to another—then one finds that $f(z)$ ends in a different place than it began, even though z itself has returned precisely to its own starting point. The range contour remains open even though the domain contour is closed.

In complex analysis, a branch point may be thought of informally as a point z_o at which a “multiple-valued function” changes values when one winds once around z_o . [46, “Branch point,” 18:10, 16 May 2006]

An analytic function like $g(z) = \sqrt{z}$ having a branch point evidently is not single-valued. It is multiple-valued. For a single z more than one distinct $g(z)$ is possible.

An analytic function like $h(z) = 1/z$, by contrast, is single-valued even though it has a pole. This function does not suffer the syndrome described. When a domain contour encircles a pole, the corresponding range contour is properly closed. Poles do not cause their functions to be multiple-valued and thus are not branch points.

Evidently $f(z) \equiv (z - z_o)^a$ has a branch point at $z = z_o$ if and only if a is not an integer. If $f(z)$ does have a branch point—if a is not an integer—then the mathematician must draw a distinction between $z_1 = z_o + \rho e^{i\phi}$ and $z_2 = z_o + \rho e^{i(\phi+2\pi)}$, *even though the two are exactly the same number*. Indeed $z_1 = z_2$, but paradoxically $f(z_1) \neq f(z_2)$.

This is difficult. It is confusing, too, until one realizes that the fact of a branch point says nothing whatsoever about the argument z . As far as z is concerned, there really is no distinction between $z_1 = z_o + \rho e^{i\phi}$ and $z_2 = z_o + \rho e^{i(\phi+2\pi)}$ —none at all. What draws the distinction is the multiple-valued function $f(z)$ which uses the argument.

It is as though I had a mad colleague who called me Thaddeus Black, until one day I happened to walk past behind his desk (rather than in front as I usually did), whereupon for some reason he began calling me Gorbag Pfufnik. I had not changed at all, but now the colleague calls me by a different name. The change isn't really in me, is it? It's in my colleague, who seems to suffer a branch point. If it is important to me to be sure that my colleague really is addressing me when he cries, “Pfufnik!” then I had better keep a running count of how many times I have turned about his desk, hadn't I, even though the number of turns is personally of no import to me.

The usual analysis strategy when one encounters a branch point is simply to avoid the point. Where an integral follows a closed contour as in § 8.8,

the strategy is to compose the contour to exclude the branch point, to shut it out. Such a strategy of avoidance usually prospers.¹⁰

8.6 Entire and meromorphic functions

Though an applied mathematician is unwise to let abstract definitions enthrall his thinking, pure mathematics nevertheless brings some technical definitions the applied mathematician can use. Two such are the definitions of *entire* and *meromorphic* functions.¹¹

A function $f(z)$ which is analytic for all finite z is an *entire function*. Examples include $f(z) = z^2$ and $f(z) = \exp z$, but not $f(z) = 1/z$ which has a pole at $z = 0$.

A function $f(z)$ which is analytic for all finite z except at isolated poles (which can be n -fold poles if n is a finite integer), which has no branch points, of which no circle of finite radius in the Argand domain plane encompasses an infinity of poles, is a *meromorphic function*. Examples include $f(z) = 1/z$, $f(z) = 1/(z+2) + 1/(z-1)^3 + 2z^2$ and $f(z) = \tan z$ —the last of which has an infinite number of poles, but of which the poles nowhere cluster in infinite numbers. The function $f(z) = \tan(1/z)$ is not meromorphic since it has an infinite number of poles within the Argand unit circle. Even the function $f(z) = \exp(1/z)$ is not meromorphic: it has only the one, isolated nonanalytic point at $z = 0$, and that point is no branch point; but the point is an *essential singularity*, having the character of an infinitifold (∞ -fold) pole.¹²

If it seems unclear that the singularities of $\tan z$ are actual poles, incidentally, then consider that

$$\tan z = \frac{\sin z}{\cos z} = -\frac{\cos w}{\sin w},$$

wherein we have changed the variable

$$w \leftarrow z - (2n+1)\frac{2\pi}{4}.$$

¹⁰Traditionally associated with branch points in complex variable theory are the notions of *branch cuts* and *Riemann sheets*. These ideas are interesting, but are not central to the analysis as developed in this book and are not covered here. The interested reader might consult a book on complex variables or advanced calculus like [21], among many others.

¹¹[45]

¹²[25]

Section 8.9 and its Table 8.1, below, give Taylor series for $\cos z$ and $\sin z$, with which

$$\tan z = \frac{-1 + w^2/2 - w^4/0x18 - \dots}{w - w^3/6 + w^5/0x78 - \dots}.$$

By long division,

$$\tan z = -\frac{1}{w} + \frac{w/3 - w^3/0x1E + \dots}{1 - w^2/6 + w^4/0x78 - \dots}.$$

(On the other hand, if it is unclear that $z = [2n + 1][2\pi/4]$ are the only singularities $\tan z$ has—that it has no singularities of which $\Im[z] \neq 0$ —then consider that the singularities of $\tan z$ occur where $\cos z = 0$, which by Euler’s formula, eqn. 5.17, occurs where $\exp[+iz] = \exp[-iz]$. This in turn is possible only if $|\exp[+iz]| = |\exp[-iz]|$, which happens only for real z .)

Sections 8.14, 8.15 and 9.6 speak further of the matter.

8.7 Extrema over a complex domain

If a function $f(z)$ is expanded by (8.19) or by other means about an analytic expansion point $z = z_o$ such that

$$f(z) = f(z_o) + \sum_{k=1}^{\infty} (a_k)(z - z_o)^k;$$

and if

$$\begin{aligned} a_k &= 0 \quad \text{for } 0 < k < K < \infty, \text{ but} \\ a_K &\neq 0, \end{aligned}$$

such that a_K , $0 < K < \infty$, is the series’ first nonzero coefficient; then, in the immediate neighborhood of the expansion point,

$$f(z) \approx f(z_o) + (a_K)(z - z_o)^K, \quad |z - z_o| \ll 1.$$

Changing $\rho'e^{i\phi'} \leftarrow z - z_o$, this is

$$f(z) \approx f(z_o) + a_K \rho'^K e^{iK\phi'}, \quad 0 \leq \rho' \ll 1. \quad (8.20)$$

Evidently one can shift the output of an analytic function $f(z)$ slightly in any desired Argand direction by shifting slightly the function’s input z . Specifically according to (8.20), to shift $f(z)$ by $\Delta f \approx \epsilon e^{i\psi}$, one can shift z

by $\Delta z \approx (\epsilon/a_K)^{1/K} e^{i(\psi+n2\pi)/K}$. Except at a nonanalytic point of $f(z)$ or in the trivial case that $f(z)$ were everywhere constant, this always works—even where $[df/dz]_{z=z_o} = 0$.

That one can shift an analytic function's output smoothly in any Argand direction whatsoever has the significant consequence that neither the real nor the imaginary part of the function—nor for that matter any linear combination $\Re[e^{-i\omega} f(z)]$ of the real and imaginary parts—can have an extremum within the interior of a domain over which the function is fully analytic. That is, *a function's extrema over a bounded analytic domain never lie within the domain's interior but always on its boundary.*^{13,14}

8.8 Cauchy's integral formula

In § 7.6 we considered the problem of vector contour integration, in which the sum value of an integration depends not only on the integration's endpoints but also on the path, or *contour*, over which the integration is done, as in Fig. 7.8. Because real scalars are confined to a single line, no alternate choice of path is possible where the variable of integration is a real scalar, so the contour problem does not arise in that case. It does however arise where the variable of integration is a *complex* scalar, because there again different paths are possible. Refer to the Argand plane of Fig. 2.5.

Consider the integral

$$S_n = \int_{z_1}^{z_2} z^{n-1} dz. \quad (8.21)$$

If z were always a real number, then by the antiderivative (§ 7.2) this integral would evaluate to $(z_2^n - z_1^n)/n$; or, in the case of $n = 0$, to $\ln(z_2/z_1)$. Inasmuch as z is complex, however, the correct evaluation is less obvious. To evaluate the integral sensibly in the latter case, one must consider some specific path of integration in the Argand plane. One must also consider the meaning of the symbol dz .

¹³Professional mathematicians tend to define the domain and its boundary more carefully.

¹⁴[7][25]

8.8.1 The meaning of the symbol dz

The symbol dz represents an infinitesimal step in some direction in the Argand plane:

$$\begin{aligned}
 dz &= [z + dz] - [z] \\
 &= [(\rho + d\rho)e^{i(\phi+d\phi)}] - [\rho e^{i\phi}] \\
 &= [(\rho + d\rho)e^{i d\phi} e^{i\phi}] - [\rho e^{i\phi}] \\
 &= [(\rho + d\rho)(1 + i d\phi)e^{i\phi}] - [\rho e^{i\phi}].
 \end{aligned}$$

Since the product of two infinitesimals is negligible even on infinitesimal scale, we can drop the $d\rho d\phi$ term.¹⁵ After canceling finite terms, we are left with the peculiar but excellent formula

$$dz = (d\rho + i\rho d\phi)e^{i\phi}. \quad (8.22)$$

8.8.2 Integrating along the contour

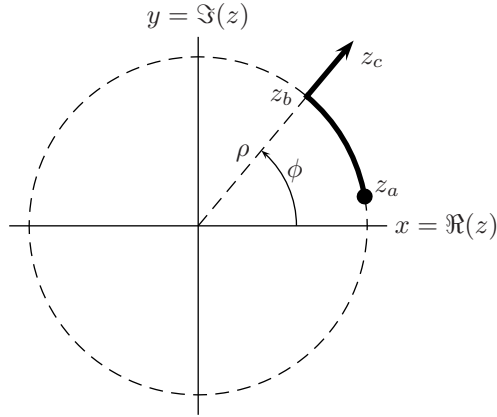
Now consider the integration (8.21) along the contour of Fig. 8.1. Integrat-

¹⁵The dropping of second-order infinitesimals like $d\rho d\phi$, added to first order infinitesimals like $d\rho$, is a standard calculus technique. One cannot *always* drop them, however. Occasionally one encounters a sum in which not only do the finite terms cancel, but also the first-order infinitesimals. In such a case, the second-order infinitesimals dominate and cannot be dropped. An example of the type is

$$\lim_{\epsilon \rightarrow 0} \frac{(1 - \epsilon)^3 + 3(1 + \epsilon) - 4}{\epsilon^2} = \lim_{\epsilon \rightarrow 0} \frac{(1 - 3\epsilon + 3\epsilon^2) + (3 + 3\epsilon) - 4}{\epsilon^2} = 3.$$

One typically notices that such a case has arisen when the dropping of second-order infinitesimals has left an ambiguous 0/0. To fix the problem, you simply go back to the step where you dropped the infinitesimal and you restore it, then you proceed from there. Otherwise there isn't much point in carrying second-order infinitesimals around. In the relatively uncommon event that you need them, you'll know it. The math itself will tell you.

Figure 8.1: A contour of integration in the Argand plane, in two segments: constant- ρ (z_a to z_b); and constant- ϕ (z_b to z_c).



ing along the constant- ϕ segment,

$$\begin{aligned}
 \int_{z_b}^{z_c} z^{n-1} dz &= \int_{\rho_b}^{\rho_c} (\rho e^{i\phi})^{n-1} (d\rho + i\rho d\phi) e^{i\phi} \\
 &= \int_{\rho_b}^{\rho_c} (\rho e^{i\phi})^{n-1} (d\rho) e^{i\phi} \\
 &= e^{in\phi} \int_{\rho_b}^{\rho_c} \rho^{n-1} d\rho \\
 &= \frac{e^{in\phi}}{n} (\rho_c^n - \rho_b^n) \\
 &= \frac{z_c^n - z_b^n}{n}.
 \end{aligned}$$

Integrating along the constant- ρ arc,

$$\begin{aligned}
 \int_{z_a}^{z_b} z^{n-1} dz &= \int_{\phi_a}^{\phi_b} (\rho e^{i\phi})^{n-1} (d\rho + i\rho d\phi) e^{i\phi} \\
 &= \int_{\phi_a}^{\phi_b} (\rho e^{i\phi})^{n-1} (i\rho d\phi) e^{i\phi} \\
 &= i\rho^n \int_{\phi_a}^{\phi_b} e^{in\phi} d\phi \\
 &= \frac{i\rho^n}{in} (e^{in\phi_b} - e^{in\phi_a}) \\
 &= \frac{z_b^n - z_a^n}{n}.
 \end{aligned}$$

Adding the two, we have that

$$\int_{z_a}^{z_c} z^{n-1} dz = \frac{z_c^n - z_a^n}{n},$$

surprisingly the same as for real z . Since any path of integration between any two complex numbers z_1 and z_2 is approximated arbitrarily closely by a succession of short constant- ρ and constant- ϕ segments, it follows generally that

$$\int_{z_1}^{z_2} z^{n-1} dz = \frac{z_2^n - z_1^n}{n}, \quad n \neq 0. \quad (8.23)$$

The applied mathematician might reasonably ask, “Was (8.23) really worth the trouble? We knew *that* already. It’s the same as for real numbers.”

Well, we really didn’t know it before deriving it, but the point is well taken nevertheless. However, notice the exemption of $n = 0$. Equation (8.23) does not hold in that case. Consider the $n = 0$ integral

$$S_0 = \int_{z_1}^{z_2} \frac{dz}{z}.$$

Following the same steps as before and using (5.7) and (2.39), we find that

$$\int_{\rho_1}^{\rho_2} \frac{dz}{z} = \int_{\rho_1}^{\rho_2} \frac{(d\rho + i\rho d\phi)e^{i\phi}}{\rho e^{i\phi}} = \int_{\rho_1}^{\rho_2} \frac{d\rho}{\rho} = \ln \frac{\rho_2}{\rho_1}. \quad (8.24)$$

This is always real-valued, but otherwise it brings no surprise. However,

$$\int_{\phi_1}^{\phi_2} \frac{dz}{z} = \int_{\phi_1}^{\phi_2} \frac{(d\rho + i\rho d\phi)e^{i\phi}}{\rho e^{i\phi}} = i \int_{\phi_1}^{\phi_2} d\phi = i(\phi_2 - \phi_1). \quad (8.25)$$

The odd thing about this is in what happens when the contour closes a complete loop in the Argand plane about the $z = 0$ pole. In this case, $\phi_2 = \phi_1 + 2\pi$, thus

$$S_0 = i2\pi,$$

even though the integration ends where it begins.

Generalizing, we have that

$$\begin{aligned} \oint (z - z_o)^{n-1} dz &= 0, \quad n \neq 0; \\ \oint \frac{dz}{z - z_o} &= i2\pi; \end{aligned} \tag{8.26}$$

where as in § 7.6 the symbol \oint represents integration about a closed contour that ends where it begins, and where it is implied that the contour loops positively (counterclockwise, in the direction of increasing ϕ) exactly once about the $z = z_o$ pole.

Notice that the formula's $i2\pi$ does not depend on the precise path of integration, but only on the fact that the path loops once positively about the pole. Notice also that nothing in the derivation of (8.23) actually requires that n be an integer, so one can write

$$\int_{z_1}^{z_2} z^{a-1} dz = \frac{z_2^a - z_1^a}{a}, \quad a \neq 0. \tag{8.27}$$

However, (8.26) does not hold in the latter case; its integral comes to zero for nonintegral a only if the contour does not enclose the branch point at $z = z_o$.

For a closed contour *which encloses no pole or other nonanalytic point*, (8.27) has that $\oint z^{a-1} dz = 0$, or with the change of variable $z - z_o \leftarrow z$,

$$\oint (z - z_o)^{a-1} dz = 0.$$

But because any analytic function can be expanded in the form $f(z) = \sum_k (c_k)(z - z_o)^{a_k-1}$ (which is just a Taylor series if the a_k happen to be positive integers), this means that

$$\oint f(z) dz = 0 \tag{8.28}$$

if $f(z)$ is everywhere analytic within the contour.¹⁶

¹⁶The careful reader will observe that (8.28)'s derivation does not explicitly handle

8.8.3 The formula

The combination of (8.26) and (8.28) is powerful. Consider the closed contour integral

$$\oint \frac{f(z)}{z - z_o} dz,$$

where the contour encloses no nonanalytic point of $f(z)$ itself but does enclose the pole of $f(z)/(z - z_o)$ at $z = z_o$. If the contour were a tiny circle of infinitesimal radius about the pole, then the integrand would reduce to $f(z_o)/(z - z_o)$; and then per (8.26),

$$\oint \frac{f(z)}{z - z_o} dz = i2\pi f(z_o). \quad (8.29)$$

But if the contour were not an infinitesimal circle but rather the larger contour of Fig. 8.2? In this case, if the dashed detour which excludes the pole is taken, then according to (8.28) the resulting integral totals zero; but the two straight integral segments evidently cancel; and similarly as we have just reasoned, the *reverse-directed* integral about the tiny detour circle is $-i2\pi f(z_o)$; so to bring the total integral to zero the integral about the main contour must be $i2\pi f(z_o)$. Thus, (8.29) holds for any positively-directed contour which once encloses a pole and no other nonanalytic point, whether the contour be small or large. Equation (8.29) is *Cauchy's integral formula*.

If the contour encloses multiple poles (§§ 2.11 and 9.6.2), then by the principle of linear superposition (§ 7.3.3),

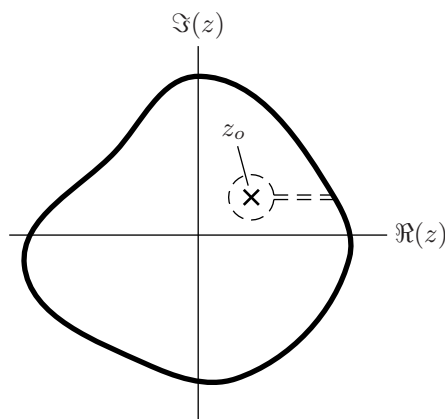
$$\oint \left[f_o(z) + \sum_k \frac{f_k(z)}{z - z_k} \right] dz = i2\pi \sum_k f_k(z_k), \quad (8.30)$$

where the $f_o(z)$ is a *regular part*;¹⁷ and again, where neither $f_o(z)$ nor any of the several $f_k(z)$ has a pole or other nonanalytic point within (or on) the

an $f(z)$ represented by a Taylor series with an infinite number of terms and a finite convergence domain (for example, $f[z] = \ln[1 - z]$). However, by § 8.2 one can transpose such a series from z_o to an overlapping convergence domain about z_1 . Let the contour's interior be divided into several cells, each of which is small enough to enjoy a single convergence domain. Integrate about each cell. Because the cells share boundaries within the contour's interior, each interior boundary is integrated twice, once in each direction, canceling. The original contour—each piece of which is an exterior boundary of some cell—is integrated once piecewise. This is the basis on which a more rigorous proof is constructed.

¹⁷[28, § 1.1]

Figure 8.2: A Cauchy contour integral.



contour. The values $f_k(z_k)$, which represent the strengths of the poles, are called *residues*. In words, (8.30) says that an integral about a closed contour in the Argand plane comes to $i2\pi$ times the sum of the residues of the poles (if any) thus enclosed. (Note however that eqn. 8.30 does not handle branch points. If there is a branch point, the contour must exclude it or the formula will not work.)

As we shall see in § 9.5, whether in the form of (8.29) or of (8.30) Cauchy's integral formula is an extremely useful result.¹⁸

8.8.4 Enclosing a multiple pole

When a complex contour of integration encloses a double, triple or other n -fold pole, the integration can be written

$$S = \oint \frac{f(z)}{(z - z_o)^{m+1}} dz,$$

where $m + 1 = n$. Expanding $f(z)$ in a Taylor series (8.19) about $z = z_o$,

$$S = \oint \sum_{k=0}^{\infty} \left(\frac{d^k f}{dz^k} \Big|_{z=z_o} \right) \frac{dz}{(k!)(z - z_o)^{m-k+1}}.$$

¹⁸[21, § 10.6][39][46, "Cauchy's integral formula," 14:13, 20 April 2006]

But according to (8.26), only the $k = m$ term contributes, so

$$\begin{aligned} S &= \oint \left(\frac{d^m f}{dz^m} \Big|_{z=z_o} \right) \frac{dz}{(m!)(z - z_o)} \\ &= \frac{1}{m!} \left(\frac{d^m f}{dz^m} \Big|_{z=z_o} \right) \oint \frac{dz}{(z - z_o)} \\ &= \frac{i2\pi}{m!} \left(\frac{d^m f}{dz^m} \Big|_{z=z_o} \right), \end{aligned}$$

where the integral is evaluated in the last step according to (8.29). Altogether,

$$\oint \frac{f(z)}{(z - z_o)^{m+1}} dz = \frac{i2\pi}{m!} \left(\frac{d^m f}{dz^m} \Big|_{z=z_o} \right), \quad m \geq 0. \quad (8.31)$$

Equation (8.31) evaluates a contour integral about an n -fold pole as (8.29) does about a single pole. (When $m = 0$, the two equations are the same.)¹⁹

8.9 Taylor series for specific functions

With the general Taylor series formula (8.19), the derivatives of Tables 5.2 and 5.3, and the observation from (4.21) that

$$\frac{d(z^a)}{dz} = az^{a-1},$$

one can calculate Taylor series for many functions. For instance, expanding about $z = 1$,

$$\begin{aligned} \ln z \Big|_{z=1} &= & \ln z \Big|_{z=1} &= & 0, \\ \frac{d}{dz} \ln z \Big|_{z=1} &= & \frac{1}{z} \Big|_{z=1} &= & 1, \\ \frac{d^2}{dz^2} \ln z \Big|_{z=1} &= & \frac{-1}{z^2} \Big|_{z=1} &= & -1, \\ \frac{d^3}{dz^3} \ln z \Big|_{z=1} &= & \frac{2}{z^3} \Big|_{z=1} &= & 2, \\ & & \vdots & & \\ \frac{d^k}{dz^k} \ln z \Big|_{z=1} &= & \frac{-(-)^k(k-1)!}{z^k} \Big|_{z=1} &= & -(-)^k(k-1)!, \quad k > 0. \end{aligned}$$

¹⁹[25][39]

With these derivatives, the Taylor series about $z = 1$ is

$$\ln z = \sum_{k=1}^{\infty} \left[-(-)^k (k-1)! \right] \frac{(z-1)^k}{k!} = - \sum_{k=1}^{\infty} \frac{(1-z)^k}{k},$$

evidently convergent for $|1-z| < 1$. (And if z lies outside the convergence domain? Several strategies are then possible. One can expand the Taylor series about a different point; but cleverer and easier is to take advantage of some convenient relationship like $\ln w = -\ln[1/w]$. Section 8.10.4 elaborates.) Using such Taylor series, one can relatively efficiently calculate actual numerical values for $\ln z$ and many other functions.

Table 8.1 lists Taylor series for a few functions of interest. All the series converge for $|z| < 1$. The $\exp z$, $\sin z$ and $\cos z$ series converge for all complex z . Among the several series, the series for $\arctan z$ is computed indirectly²⁰ by way of Table 5.3 and (2.33):

$$\begin{aligned} \arctan z &= \int_0^z \frac{1}{1+w^2} dw \\ &= \int_0^z \sum_{k=0}^{\infty} (-)^k w^{2k} dw \\ &= \sum_{k=0}^{\infty} \frac{(-)^k z^{2k+1}}{2k+1}. \end{aligned}$$

8.10 Error bounds

One naturally cannot actually sum a Taylor series to an infinite number of terms. One must add some finite number of terms, then quit—which raises the question: how many terms are enough? How can one know that one has added adequately many terms; that the remaining terms, which constitute the tail of the series, are sufficiently insignificant? How can one set error bounds on the truncated sum?

8.10.1 Examples

Some series alternate sign. For these it is easy if the numbers involved happen to be real. For example, from Table 8.1,

$$\ln \frac{3}{2} = \ln \left(1 + \frac{1}{2} \right) = \frac{1}{(1)(2^1)} - \frac{1}{(2)(2^2)} + \frac{1}{(3)(2^3)} - \frac{1}{(4)(2^4)} + \cdots$$

²⁰[37, § 11-7]

Table 8.1: Taylor series.

$$\begin{aligned}
f(z) &= \sum_{k=0}^{\infty} \left(\left. \frac{d^k f}{dz^k} \right|_{z=z_o} \right) \prod_{j=1}^k \frac{z - z_o}{j} \\
(1+z)^{a-1} &= \sum_{k=0}^{\infty} \prod_{j=1}^k \left(\frac{a}{j} - 1 \right) z \\
\exp z &= \sum_{k=0}^{\infty} \prod_{j=1}^k \frac{z}{j} = \sum_{k=0}^{\infty} \frac{z^k}{k!} \\
\sin z &= \sum_{k=0}^{\infty} \left[z \prod_{j=1}^k \frac{-z^2}{(2j)(2j+1)} \right] \\
\cos z &= \sum_{k=0}^{\infty} \prod_{j=1}^k \frac{-z^2}{(2j-1)(2j)} \\
\sinh z &= \sum_{k=0}^{\infty} \left[z \prod_{j=1}^k \frac{z^2}{(2j)(2j+1)} \right] \\
\cosh z &= \sum_{k=0}^{\infty} \prod_{j=1}^k \frac{z^2}{(2j-1)(2j)} \\
-\ln(1-z) &= \sum_{k=1}^{\infty} \frac{1}{k} \prod_{j=1}^k z = \sum_{k=1}^{\infty} \frac{z^k}{k} \\
\arctan z &= \sum_{k=0}^{\infty} \frac{1}{2k+1} \left[z \prod_{j=1}^k (-z^2) \right] = \sum_{k=0}^{\infty} \frac{(-)^k z^{2k+1}}{2k+1}
\end{aligned}$$

Each term is smaller in magnitude than the last, so the true value of $\ln(3/2)$ necessarily lies between the sum of the series to n terms and the sum to $n + 1$ terms. The last and next partial sums bound the result. After three terms, for instance,

$$S_3 - \frac{1}{(4)(2^4)} < \ln \frac{3}{2} < S_3,$$

$$S_3 = \frac{1}{(1)(2^1)} - \frac{1}{(2)(2^2)} + \frac{1}{(3)(2^3)}.$$

Other series however do not alternate sign. For example,

$$\begin{aligned} \ln 2 &= -\ln \frac{1}{2} = -\ln \left(1 - \frac{1}{2}\right) = S_4 + R_4, \\ S_4 &= \frac{1}{(1)(2^1)} + \frac{1}{(2)(2^2)} + \frac{1}{(3)(2^3)} + \frac{1}{(4)(2^4)}, \\ R_4 &= \frac{1}{(5)(2^5)} + \frac{1}{(6)(2^6)} + \cdots \end{aligned}$$

The basic technique in such a case is to find a replacement series (or integral) R'_n which one can collapse analytically, each of whose terms equals or exceeds in magnitude the corresponding term of R_n . For the example, one might choose

$$R'_4 = \frac{1}{5} \sum_{k=5}^{\infty} \frac{1}{2^k} = \frac{2}{(5)(2^5)},$$

wherein (2.34) had been used to collapse the summation. Then,

$$S_4 < \ln 2 < S_4 + R'_4.$$

For real $0 \leq x < 1$ generally,

$$\begin{aligned} S_n &< -\ln(1-x) < S_n + R'_n, \\ S_n &\equiv \sum_{k=1}^n \frac{x^k}{k}, \\ R'_n &\equiv \sum_{k=n+1}^{\infty} \frac{x^k}{n+1} = \frac{x^{n+1}}{(n+1)(1-x)}. \end{aligned}$$

Many variations and refinements are possible, some of which we shall meet in the rest of the section, but that is the basic technique: to add several

terms of the series to establish a lower bound, then to overestimate the remainder of the series to establish an upper bound. The overestimate R'_n *majorizes* the series' true remainder R_n . Notice that the R'_n in the example is a fairly small number, and that it would have been a lot smaller yet had we included a few more terms in S_n (for instance, $n = 0x40$ would have bound $\ln 2$ tighter than the limit of a computer's typical **double**-type floating-point accuracy). The technique usually works well in practice for this reason.

8.10.2 Majorization

To *majorize* in mathematics is to be, or to replace by virtue of being, everywhere at least as great as. This is best explained by example. Consider the summation

$$S = \sum_{k=1}^{\infty} \frac{1}{k^2} = 1 + \frac{1}{2^2} + \frac{1}{3^2} + \frac{1}{4^2} + \cdots$$

The exact value this summation totals to is unknown to us, but the summation does rather resemble the integral (refer to Table 7.1)

$$I = \int_1^{\infty} \frac{dx}{x^2} = -\frac{1}{x} \Big|_1^{\infty} = 1.$$

Figure 8.3 plots S and I together as areas—or more precisely, plots $S - 1$ and I together as areas (the summation's first term is omitted). As the plot shows, the unknown area $S - 1$ cannot possibly be as great as the known area I . In symbols, $S - 1 < I = 1$; or,

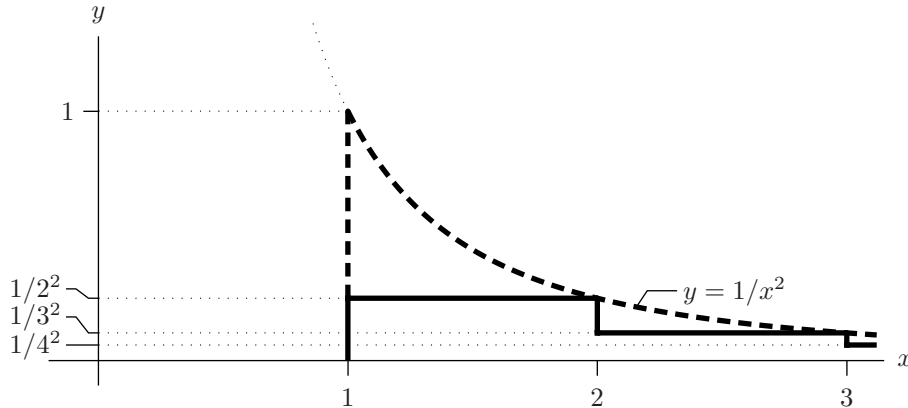
$$S < 2.$$

The integral I majorizes the summation $S - 1$, thus guaranteeing the absolute upper limit on S . (Of course $S < 2$ is a very loose limit, but that isn't the point of the example. In practical calculation, one would let a computer add many terms of the series first numerically, and only then majorize the remainder. Even so, cleverer ways to majorize the remainder of this particular series will occur to the reader, such as in representing the terms graphically—not as flat-topped rectangles—but as slant-topped trapezoids, shifted in the figure a half unit rightward.)

Majorization serves surely to bound an unknown quantity by a larger, known quantity. Reflecting, *minorization*²¹ serves surely to bound an unknown quantity by a smaller, known quantity. The quantities in question

²¹The author does not remember ever encountering the word *minorization* heretofore in print, but as a reflection of *majorization* the word seems logical. This book at least will use the word where needed. You can use it too if you like.

Figure 8.3: Majorization. The area I between the dashed curve and the x axis majorizes the area $S - 1$ between the staircase curve and the x axis, because the height of the dashed curve is everywhere at least as great as that of the staircase curve.



are often integrals and/or series summations, the two of which are akin as Fig. 8.3 illustrates. The choice of whether to majorize a particular unknown quantity by an integral or by a series summation depends on the convenience of the problem at hand.

The series S of this subsection is interesting, incidentally. It is a *harmonic series* rather than a power series, because although its terms do decrease in magnitude it has no z^k factor (or seen from another point of view, it does have a z^k factor, but $z = 1$), and the ratio of adjacent terms' magnitudes approaches unity as k grows. Harmonic series can be hard to sum accurately, but clever majorization can help (and if majorization does not help enough, the series transformations of Ch. [not yet written] can help even more).

8.10.3 Geometric majorization

Harmonic series can be hard to sum as § 8.10.3 has observed; but more common than harmonic series are true power series, which are easier to sum in that they include a z^k factor. There is no one, ideal bound which works equally well for all power series. However, the point of establishing a bound is not to sum a power series exactly but rather to fence the sum within some

sufficiently (rather than optimally) small neighborhood. A simple, general bound which works quite adequately for most power series encountered in practice, including among many others all the Taylor series of Table 8.1, is the *geometric majorization*

$$|\epsilon_n| < \frac{|\tau_{n+1}|}{1 - |\rho|}. \quad (8.32)$$

Here, τ_k represents the power series' k th-order term (in Table 8.1's series for $\exp z$, for example, $\tau_k = z^k/[k!]$). The $|\rho|$ is a positive real number chosen, preferably as small as possible, such that

$$|\rho| \geq \left| \frac{\tau_{k+1}}{\tau_k} \right| \quad \text{for all } k > n, \quad (8.33)$$

$$\begin{aligned} |\rho| &> \left| \frac{\tau_{k+1}}{\tau_k} \right| \quad \text{for at least one } k > n, \\ |\rho| &< 1; \end{aligned} \quad (8.34)$$

which is to say, more or less, such that each term in the series' tail is smaller than the last by at least a factor of $|\rho|$. Given these definitions, if

$$\begin{aligned} S_n &\equiv \sum_{k=0}^n \tau_k, \\ \epsilon_n &\equiv S_n - S_\infty, \end{aligned} \quad (8.35)$$

where S_∞ represents the true, exact (but uncalculatable, unknown) infinite series sum, then (2.34) and (3.20) imply the geometric majorization (8.32).

If the last paragraph seems abstract, a pair of concrete examples should serve to clarify. First, if the Taylor series

$$-\ln(1 - z) = \sum_{k=1}^{\infty} \frac{z^k}{k}$$

of Table 8.1 is truncated after the n th-order term, then

$$\begin{aligned} -\ln(1 - z) &\approx \sum_{k=1}^n \frac{z^k}{k}, \\ |\epsilon_n| &< \frac{|z^{n+1}|/(n+1)}{1 - |z|}, \end{aligned}$$

where ϵ_n is the error in the truncated sum. Here, $|\tau_{k+1}/\tau_k| = [k/(k+1)]|z| < |z|$, so we have chosen $|\rho| = |z|$.

Second, if the Taylor series

$$\exp z = \sum_{k=0}^{\infty} \prod_{j=1}^k \frac{z}{j} = \sum_{k=0}^{\infty} \frac{z^k}{k!}$$

also of Table 8.1 is truncated after the n th-order term, where $n + 2 > |z|$, then

$$\begin{aligned} \exp z &\approx \sum_{k=0}^n \prod_{j=1}^k \frac{z}{j} = \sum_{k=0}^n \frac{z^k}{k!}, \\ |\epsilon_n| &< \frac{|z^{n+1}|/(n+1)!}{1 - |z|/(n+2)}. \end{aligned}$$

Here, $|\tau_{k+1}/\tau_k| = |z|/(k+1)$, whose maximum value for all $k > n$ occurs when $k = n + 1$, so we have chosen $|\rho| = |z|/(n+2)$.

8.10.4 Calculation outside the fast convergence domain

Used directly, the Taylor series of Table 8.1 tend to converge slowly for some values of z and not at all for others. The series for $-\ln(1-z)$ and $(1+z)^{a-1}$ for instance each converge for $|z| < 1$ (though slowly for $|z| \approx 1$); whereas each series diverges when asked to compute a quantity like $-\ln 3$ or 3^{a-1} directly. To shift the series' expansion points per § 8.2 is one way to seek convergence, but for nonentire functions (§ 8.6) like these a more probably profitable strategy is to find and exploit some property of the functions to transform their arguments, such as

$$\begin{aligned} -\ln \gamma &= \ln \frac{1}{\gamma}, \\ \gamma^{a-1} &= \frac{1}{(1/\gamma)^{a-1}}, \end{aligned}$$

which leave the respective Taylor series to compute quantities like $-\ln(1/3)$ and $(1/3)^{a-1}$ they can handle.

Let $f(1+\zeta)$ be a function whose Taylor series about $\zeta = 0$ converges for

$|\zeta| < 1$ and which obeys properties of the forms²²

$$\begin{aligned} f(\gamma) &= g \left[f \left(\frac{1}{\gamma} \right) \right], \\ f(\alpha\gamma) &= h[f(\alpha), f(\gamma)], \end{aligned} \quad (8.36)$$

where $g[\cdot]$ and $h[\cdot, \cdot]$ are functions we know how to compute like $g[\cdot] = -[\cdot]$ or $g[\cdot] = 1/[\cdot]$; and like $h[\cdot, \cdot] = [\cdot] + [\cdot]$ or $h[\cdot, \cdot] = [\cdot][\cdot]$. Identifying

$$\begin{aligned} \frac{1}{\gamma} &= 1 + \zeta, \\ \frac{1 - \gamma}{\gamma} &= \zeta, \end{aligned} \quad (8.37)$$

we have that

$$f(\gamma) = g \left[f \left(1 + \frac{1 - \gamma}{\gamma} \right) \right], \quad (8.38)$$

whose convergence domain $|\zeta| < 1$ is $|1 - \gamma|/|\gamma| < 1$, which is $|\gamma - 1| < |\gamma|$ or in other words

$$\Re(\gamma) > \frac{1}{2}.$$

Although the transformation from ζ to γ has not lifted the convergence limit altogether, it has apparently opened it to a broader domain.

Though this writer knows no way to lift the convergence limit altogether that does not cause more problems than it solves, one can take advantage of the $h[\cdot, \cdot]$ property of (8.36) to sidestep the limit, computing $f(\omega)$ indirectly for any $\omega \neq 0$ by any of several tactics. One nonoptimal but entirely effective tactic is represented by the equations

$$\begin{aligned} \omega &\equiv i^n 2^m \gamma, \\ |\Im(\gamma)| &\leq \Re(\gamma), \\ 1 &\leq \Re(\gamma) < 2, \end{aligned} \quad (8.39)$$

whereupon the formula²³

$$f(\omega) = h[f(i^n 2^m), f(\gamma)] \quad (8.40)$$

²²This paragraph's notation is necessarily abstract. To make it seem more concrete, consider that the function $f(1 + \zeta) = -\ln(1 - z)$ has $z = -\zeta$, $f(\gamma) = g[f(1/\gamma)] = -f(1/\gamma)$ and $f(\alpha\gamma) = h[f(\alpha), f(\gamma)] = f(\alpha) + f(\gamma)$; and that the function $f(1 + \zeta) = (1 + z)^{a-1}$ has $z = \zeta$, $f(\gamma) = g[f(1/\gamma)] = 1/f(1/\gamma)$ and $f(\alpha\gamma) = h[f(\alpha), f(\gamma)] = f(\alpha)f(\gamma)$.

²³Equation (8.40) leaves open the question of how to compute $f(i^n 2^m)$, but for the example functions at least this is not hard. For the logarithm, $-\ln(i^n 2^m) = m \ln(1/2) - in(2\pi/4)$. For the power, $(i^n 2^m)^{a-1} = \text{cis}[(n2\pi/4)(a-1)]/[(1/2)^{a-1}]^m$. The sine and cosine in the cis function are each calculated directly by Taylor series, as are the numbers $\ln(1/2)$ and $(1/2)^{a-1}$. The number 2π , we have not calculated yet, but shall in § 8.11.

calculates $f(\omega)$ fast for any $\omega \neq 0$.

Notice how (8.39) fences γ within a comfortable zone, keeping γ moderately small in magnitude but never too near the $\Re(\gamma) = 1/2$ frontier in the Argand plane. In theory all finite γ rightward of the frontier let the Taylor series converge, but extreme γ of any kind let the series converge only slowly (and due to compound floating-point rounding error inaccurately) inasmuch as they imply that $|\zeta| \approx 1$. Besides allowing all $\omega \neq 0$, the tactic (8.39) also thus significantly speeds series convergence.

The method and tactic of (8.36) through (8.40) are useful in themselves and also illustrative generally. Of course, most nonentire functions lack properties of the specific kinds that (8.36) demands, but such functions may have other properties one can analogously exploit.²⁴

8.10.5 Nonconvergent series

Variants of this section's techniques can be used to prove that a series does not converge at all. For example,

$$\sum_{k=1}^{\infty} \frac{1}{k}$$

does not converge because

$$\frac{1}{k} > \int_k^{k+1} \frac{d\tau}{\tau};$$

hence,

$$\sum_{k=1}^{\infty} \frac{1}{k} > \sum_{k=1}^{\infty} \int_k^{k+1} \frac{d\tau}{\tau} = \int_1^{\infty} \frac{d\tau}{\tau} = \ln \infty.$$

²⁴To draw another example from Table 8.1, consider that

$$\begin{aligned} \arctan \omega &= \alpha + \arctan \zeta, \\ \zeta &\equiv \frac{\omega \cos \alpha - \sin \alpha}{\omega \sin \alpha + \cos \alpha}, \end{aligned}$$

where $\arctan \omega$ is interpreted as the geometrical angle the vector $\hat{\mathbf{x}} + \hat{\mathbf{y}}\omega$ makes with $\hat{\mathbf{x}}$. Axes are rotated per (3.7) through some angle α to reduce the tangent from ω to ζ , thus causing the Taylor series to converge faster or indeed to converge at all.

Any number of further examples and tactics of the kind will occur to the creative reader, shrinking a function's argument by some convenient means before feeding the argument to the function's Taylor series.

8.10.6 Remarks

The study of error bounds is not a matter of rules and formulas so much as of ideas, suggestions and tactics. There is normally no such thing as an optimal error bound—with sufficient cleverness, some tighter bound can usually be discovered—but often easier and more effective than cleverness is simply to add a few extra terms into the series before truncating it (that is, to increase n a little). To eliminate the error entirely usually demands adding an infinite number of terms, which is impossible; but since eliminating the error entirely also requires recording the sum to infinite precision, which is impossible anyway, eliminating the error entirely is not normally a goal one seeks. To eliminate the error to the 0x34-bit (sixteen-decimal place) precision of a computer's `double`-type floating-point representation typically requires something like 0x34 terms—if the series be wisely composed and if care be taken to keep z moderately small and reasonably distant from the edge of the series' convergence domain. Besides, few engineering applications really use much more than 0x10 bits (five decimal places) in any case. Perfect precision is impossible, but adequate precision is usually not hard to achieve.

Occasionally nonetheless a series arises for which even adequate precision is quite hard to achieve. An infamous example is

$$S = - \sum_{k=1}^{\infty} \frac{(-)^k}{\sqrt{k}} = 1 - \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{3}} - \frac{1}{\sqrt{4}} + \cdots,$$

which obviously converges, but sum it if you can! It is not easy to do.

Before closing the section, we ought to arrest one potential agent of terminological confusion. The “error” in a series summation's error bounds is unrelated to the error of probability theory. The English word “error” is thus overloaded here. A series sum converges to a definite value, and to the same value every time the series is summed; no chance is involved. It is just that we do not necessarily know exactly what that value is. What we can do, by this section's techniques or perhaps by other methods, is to establish a definite neighborhood in which the unknown value is sure to lie; and we can make that neighborhood as tight as we want, merely by including a sufficient number of terms in the sum.

The topic of series error bounds is what G.S. Brown refers to as “trick-based.”²⁵ There is no general answer to the error-bound problem, but there are several techniques which help, some of which this section has introduced. Other techniques, we shall meet later in the book as the need for them arises.

²⁵[5]

8.11 Calculating 2π

The Taylor series for $\arctan z$ in Table 8.1 implies a neat way of calculating the constant 2π . We already know that $\tan 2\pi/8 = 1$, or in other words that

$$\arctan 1 = \frac{2\pi}{8}.$$

Applying the Taylor series, we have that

$$2\pi = 8 \sum_{k=0}^{\infty} \frac{(-)^k}{2k+1}. \quad (8.41)$$

The series (8.41) is simple but converges extremely slowly. Much faster convergence is given by angles smaller than $2\pi/8$. For example, from Table 3.2,

$$\arctan \frac{\sqrt{3}-1}{\sqrt{3}+1} = \frac{2\pi}{0x18}.$$

Applying the Taylor series at this angle, we have that²⁶

$$2\pi = 0x18 \sum_{k=0}^{\infty} \frac{(-)^k}{2k+1} \left(\frac{\sqrt{3}-1}{\sqrt{3}+1} \right)^{2k+1} \approx 0x6.487F. \quad (8.42)$$

8.12 Odd and even functions

An *odd function* is one for which $f(-z) = -f(z)$. Any function whose Taylor series about $z_o = 0$ includes only odd-order terms is an odd function. Examples of odd functions include z^3 and $\sin z$.

An *even function* is one for which $f(-z) = f(z)$. Any function whose Taylor series about $z_o = 0$ includes only even-order terms is an even function. Examples of even functions include z^2 and $\cos z$.

Odd and even functions are interesting because of the symmetry they bring—the plot of a real-valued odd function being symmetric about a point, the plot of a real-valued even function being symmetric about a line. Many

²⁶The writer is given to understand that clever mathematicians have invented subtle, still much faster-converging iterative schemes toward 2π . However, there is fast and there is fast. The relatively straightforward series this section gives converges to the best accuracy of your computer's floating-point register within a paltry 0x40 iterations or so—and, after all, you only need to compute the numerical value of 2π once.

Admittedly, the writer supposes that useful lessons lurk in the clever mathematics underlying the subtle schemes, but such schemes are not covered here.

functions are neither odd nor even, of course, but one can always split an analytic function into two components—one odd, the other even—by the simple expedient of sorting the odd-order terms from the even-order terms in the function's Taylor series. For example, $\exp z = \sinh z + \cosh z$.

8.13 Trigonometric poles

The singularities of the trigonometric functions are single poles of residue ± 1 or $\pm i$. For the circular trigonometrics, all the poles lie along the real number line; for the hyperbolic trigonometrics, along the imaginary. Specifically, of the eight trigonometric functions

$$\frac{1}{\sin z}, \frac{1}{\cos z}, \frac{1}{\tan z}, \tan z, \\ \frac{1}{\sinh z}, \frac{1}{\cosh z}, \frac{1}{\tanh z}, \tanh z,$$

the poles and their respective residues are

$$\begin{aligned} \left. \frac{z - k\pi}{\sin z} \right|_{z = k\pi} &= (-)^k, \\ \left. \frac{z - (k - 1/2)\pi}{\cos z} \right|_{z = (k - 1/2)\pi} &= (-)^k, \\ \left. \frac{z - k\pi}{\tan z} \right|_{z = k\pi} &= 1, \\ [z - (k - 1/2)\pi] \tan z \Big|_{z = (k - 1/2)\pi} &= -1, \\ \left. \frac{z - ik\pi}{\sinh z} \right|_{z = ik\pi} &= (-)^k, \\ \left. \frac{z - i(k - 1/2)\pi}{\cosh z} \right|_{z = i(k - 1/2)\pi} &= i(-)^k, \\ \left. \frac{z - ik\pi}{\tanh z} \right|_{z = ik\pi} &= 1, \\ [z - i(k - 1/2)\pi] \tanh z \Big|_{z = i(k - 1/2)\pi} &= 1. \end{aligned} \tag{8.43}$$

To support (8.43)'s claims, we shall marshal the identities of Tables 5.1 and 5.2 plus l'Hôpital's rule (4.29). Before calculating residues and such, however, we should like to verify that the poles (8.43) lists are in fact the only poles that there are; that we have forgotten no poles. Consider for instance the function $1/\sin z = i2/(e^{iz} - e^{-iz})$. This function evidently goes

infinite only when $e^{iz} = e^{-iz}$, which is possible only for real z ; but for real z , the sine function's very definition establishes the poles $z = k\pi$ (refer to Fig. 3.1). With the observations from Table 5.1 that $i \sinh z = \sin iz$ and $\cosh z = \cos iz$, similar reasoning for each of the eight trigonometrics forbids poles other than those (8.43) lists. Satisfied that we have forgotten no poles, therefore, we finally apply l'Hôpital's rule to each of the ratios

$$\frac{z - k\pi}{\sin z}, \frac{z - (k - 1/2)\pi}{\cos z}, \frac{z - k\pi}{\tan z}, \frac{z - (k - 1/2)\pi}{1/\tan z},$$

$$\frac{z - ik\pi}{\sinh z}, \frac{z - i(k - 1/2)\pi}{\cosh z}, \frac{z - ik\pi}{\tanh z}, \frac{z - i(k - 1/2)\pi}{1/\tanh z}$$

to reach (8.43).

Trigonometric poles evidently are special only in that a trigonometric function has an infinite number of them. The poles are ordinary, single poles, with residues, subject to Cauchy's integral formula and so on. The trigonometrics are meromorphic functions (§ 8.6) for this reason.²⁷

The six simpler trigonometrics $\sin z$, $\cos z$, $\sinh z$, $\cosh z$, $\exp z$ and $\csc z$ —conspicuously excluded from this section's gang of eight—have no poles for finite z , because e^z , e^{iz} , $e^z \pm e^{-z}$ and $e^{iz} \pm e^{-iz}$ then likewise are finite. These simpler trigonometrics are not only meromorphic but also entire. Observe however that the *inverse* trigonometrics are multiple-valued and have branch points, and thus are not meromorphic at all.

8.14 The Laurent series

Any analytic function can be expanded in a Taylor series, but never about a pole or branch point of the function. Sometimes one nevertheless wants to expand at least about a pole. Consider for example expanding

$$f(z) = \frac{e^{-z}}{1 - \cos z} \tag{8.44}$$

²⁷[25]

about the function's pole at $z = 0$. Expanding dividend and divisor separately,

$$\begin{aligned}
 f(z) &= \frac{1 - z + z^2/2 - z^3/6 + \cdots}{z^2/2 - z^4/0x18 + \cdots} \\
 &= \frac{\sum_{j=0}^{\infty} [(-)^j z^j / j!]}{-\sum_{k=1}^{\infty} (-)^k z^{2k} / (2k)!} \\
 &= \frac{\sum_{k=0}^{\infty} [-z^{2k} / (2k)! + z^{2k+1} / (2k+1)!]}{\sum_{k=1}^{\infty} (-)^k z^{2k} / (2k)!}.
 \end{aligned}$$

By long division,

$$\begin{aligned}
 f(z) &= \frac{2}{z^2} - \frac{2}{z} + \left\{ \left[-\frac{2}{z^2} + \frac{2}{z} \right] \sum_{k=1}^{\infty} \frac{(-)^k z^{2k}}{(2k)!} \right. \\
 &\quad \left. + \sum_{k=0}^{\infty} \left[-\frac{z^{2k}}{(2k)!} + \frac{z^{2k+1}}{(2k+1)!} \right] \right\} / \sum_{k=1}^{\infty} \frac{(-)^k z^{2k}}{(2k)!} \\
 &= \frac{2}{z^2} - \frac{2}{z} + \left\{ \sum_{k=1}^{\infty} \left[-\frac{(-)^k 2 z^{2k-2}}{(2k)!} + \frac{(-)^k 2 z^{2k-1}}{(2k)!} \right] \right. \\
 &\quad \left. + \sum_{k=0}^{\infty} \left[-\frac{z^{2k}}{(2k)!} + \frac{z^{2k+1}}{(2k+1)!} \right] \right\} / \sum_{k=1}^{\infty} \frac{(-)^k z^{2k}}{(2k)!} \\
 &= \frac{2}{z^2} - \frac{2}{z} + \left\{ \sum_{k=0}^{\infty} \left[\frac{(-)^k 2 z^{2k}}{(2k+2)!} - \frac{(-)^k 2 z^{2k+1}}{(2k+2)!} \right] \right. \\
 &\quad \left. + \sum_{k=0}^{\infty} \left[-\frac{z^{2k}}{(2k)!} + \frac{z^{2k+1}}{(2k+1)!} \right] \right\} / \sum_{k=1}^{\infty} \frac{(-)^k z^{2k}}{(2k)!} \\
 &= \frac{2}{z^2} - \frac{2}{z} + \sum_{k=0}^{\infty} \left[\frac{-(2k+1)(2k+2) + (-)^k 2}{(2k+2)!} z^{2k} \right. \\
 &\quad \left. + \frac{(2k+2) - (-)^k 2}{(2k+2)!} z^{2k+1} \right] / \sum_{k=1}^{\infty} \frac{(-)^k z^{2k}}{(2k)!}.
 \end{aligned}$$

The remainder's $k = 0$ terms now disappear as intended; so, factoring z^2/z^2 from the division leaves

$$f(z) = \frac{2}{z^2} - \frac{2}{z} + \sum_{k=0}^{\infty} \left[\frac{(2k+3)(2k+4) + (-)^k 2}{(2k+4)!} z^{2k} - \frac{(2k+4) + (-)^k 2}{(2k+4)!} z^{2k+1} \right] \bigg/ \sum_{k=0}^{\infty} \frac{(-)^k z^{2k}}{(2k+2)!}. \quad (8.45)$$

One can continue dividing to extract further terms if desired, and if all the terms

$$f(z) = \frac{2}{z^2} - \frac{2}{z} + \frac{7}{6} - \frac{z}{2} + \cdots$$

are extracted the result is the *Laurent series* proper,

$$f(z) = \sum_{k=K}^{\infty} (a_k)(z - z_o)^k, \quad K \leq 0. \quad (8.46)$$

However for many purposes (as in eqn. 8.45) the partial Laurent series

$$f(z) = \sum_{k=K}^{-1} (a_k)(z - z_o)^k + \frac{\sum_{k=0}^{\infty} (b_k)(z - z_o)^k}{\sum_{k=0}^{\infty} (c_k)(z - z_o)^k}, \quad K \leq 0, \quad c_0 \neq 0, \quad (8.47)$$

suffices and may even be preferable. In either form,

$$f(z) = \sum_{k=K}^{-1} (a_k)(z - z_o)^k + f_o(z - z_o), \quad K \leq 0, \quad (8.48)$$

where, unlike $f(z)$, $f_o(z - z_o)$ is analytic at $z = z_o$. The $f_o(z - z_o)$ of (8.48) is $f(z)$'s *regular part* at $z = z_o$.

The ordinary Taylor series diverges at a function's pole. Handling the pole separately, the Laurent series remedies this defect.^{28,29}

Sections 9.5 and 9.6 tell more about poles generally, including multiple poles like the one in the example here.

²⁸The professional mathematician's treatment of the Laurent series usually begins by defining an annular convergence domain (a convergence domain bounded without by a large circle and within by a small) in the Argand plane. From an applied point of view however what interests us is the basic technique to remove the poles from an otherwise analytic function. Further rigor is left to the professionals.

²⁹[21, § 10.8][12, § 2.5]

8.15 Taylor series in $1/z$

A little imagination helps the Taylor series a lot. The Laurent series of § 8.14 represents one way to extend the Taylor series. Several other ways are possible. The typical trouble one has with the Taylor series is that a function's poles and branch points limit the series' convergence domain. Thinking flexibly, however, one can often evade the trouble.

Consider the function

$$f(z) = \frac{\sin(1/z)}{\cos z}.$$

This function has a nonanalytic point of a most peculiar nature at $z = 0$. The point is an essential singularity, and one cannot expand the function directly about it. One could expand the function directly about some other point like $z = 1$, but calculating the Taylor coefficients would take a lot of effort and, even then, the resulting series would suffer a straitly limited convergence domain. All that however tries too hard. Depending on the application, it may suffice to write

$$f(z) = \frac{\sin w}{\cos z}, \quad w \equiv \frac{1}{z}.$$

This is

$$f(z) = \frac{z^{-1} - z^{-3}/3! + z^{-5}/5! - \dots}{1 - z^2/2! + z^4/4! - \dots},$$

which is all one needs to calculate $f(z)$ numerically—and may be all one needs for analysis, too.

As an example of a different kind, consider

$$g(z) = \frac{1}{(z-2)^2}.$$

Most often, one needs no Taylor series to handle such a function (one simply does the indicated arithmetic). Suppose however that a Taylor series specifically about $z = 0$ were indeed needed for some reason. Then by (8.1) and (4.2),

$$g(z) = \frac{1/4}{(1-z/2)^2} = \frac{1}{4} \sum_{k=0}^{\infty} \binom{1+k}{1} \left(\frac{z}{2}\right)^k = \sum_{k=0}^{\infty} \frac{k+1}{2^{k+2}} z^k,$$

That expansion is good only for $|z| < 2$, but for $|z| > 2$ we also have that

$$g(z) = \frac{1/z^2}{(1-2/z)^2} = \frac{1}{z^2} \sum_{k=0}^{\infty} \binom{1+k}{1} \left(\frac{2}{z}\right)^k = \sum_{k=2}^{\infty} \frac{2^{k-2}(k-1)}{z^k},$$

which expands in negative rather than positive powers of z . Note that we have computed the two series for $g(z)$ without ever actually taking a derivative.

Neither of the section's examples is especially interesting in itself, but their point is that it often pays to think flexibly in extending and applying the Taylor series. One is not required immediately to take the Taylor series of a function as it presents itself; one can first change variables or otherwise rewrite the function in some convenient way, then take the Taylor series either of the whole function at once or of pieces of it separately. One can expand in negative powers of z equally validly as in positive powers. And, though taking derivatives per (8.19) may be the canonical way to determine Taylor coefficients, any effective means to find the coefficients suffices.

8.16 The multidimensional Taylor series

Equation (8.19) has given the Taylor series for functions of a single variable. The idea of the Taylor series does not differ where there are two or more independent variables, only the details are a little more complicated. For example, consider the function $f(z_1, z_2) = z_1^2 + z_1 z_2 + 2z_2$, which has terms z_1^2 and $2z_2$ —these we understand—but also has the cross-term $z_1 z_2$ for which the relevant derivative is the cross-derivative $\partial^2 f / \partial z_1 \partial z_2$. Where two or more independent variables are involved, one must account for the cross-derivatives, too.

With this idea in mind, the multidimensional Taylor series is

$$f(\mathbf{z}) = \sum_{\mathbf{k}} \left(\frac{\partial^{\mathbf{k}} f}{\partial \mathbf{z}^{\mathbf{k}}} \Big|_{\mathbf{z}=\mathbf{z}_o} \right) \frac{(\mathbf{z} - \mathbf{z}_o)^{\mathbf{k}}}{\mathbf{k}!}. \quad (8.49)$$

Well, that's neat. What does it mean?

- The \mathbf{z} is a *vector*³⁰ incorporating the several independent variables z_1, z_2, \dots, z_N .
- The \mathbf{k} is a nonnegative integer vector of N counters— k_1, k_2, \dots, k_N —one for each of the independent variables. Each of the k_n runs independently from 0 to ∞ , and every permutation is possible. For example,

³⁰In this generalized sense of the word, a *vector* is an ordered set of N elements. The geometrical vector $\mathbf{v} = \hat{\mathbf{x}}x + \hat{\mathbf{y}}y + \hat{\mathbf{z}}z$ of § 3.3, then, is a vector with $N = 3$, $v_1 = x$, $v_2 = y$ and $v_3 = z$.

if $N = 2$ then

$$\begin{aligned}
 \mathbf{k} &= (k_1, k_2) \\
 &= (0, 0), (0, 1), (0, 2), (0, 3), \dots; \\
 &\quad (1, 0), (1, 1), (1, 2), (1, 3), \dots; \\
 &\quad (2, 0), (2, 1), (2, 2), (2, 3), \dots; \\
 &\quad (3, 0), (3, 1), (3, 2), (3, 3), \dots; \\
 &\quad \dots
 \end{aligned}$$

- The $\partial^{\mathbf{k}} f / \partial \mathbf{z}^{\mathbf{k}}$ represents the \mathbf{k} th cross-derivative of $f(\mathbf{z})$, meaning that

$$\frac{\partial^{\mathbf{k}} f}{\partial \mathbf{z}^{\mathbf{k}}} \equiv \left(\prod_{n=1}^N \frac{\partial^{k_n}}{(\partial z_n)^{k_n}} \right) f.$$

- The $(\mathbf{z} - \mathbf{z}_o)^{\mathbf{k}}$ represents

$$(\mathbf{z} - \mathbf{z}_o)^{\mathbf{k}} \equiv \prod_{n=1}^N (z_n - z_{on})^{k_n}.$$

- The $\mathbf{k}!$ represents

$$\mathbf{k}! \equiv \prod_{n=1}^N k_n!.$$

With these definitions, the multidimensional Taylor series (8.49) yields all the right derivatives and cross-derivatives at the expansion point $\mathbf{z} = \mathbf{z}_o$. Thus within some convergence domain about $\mathbf{z} = \mathbf{z}_o$, the multidimensional Taylor series (8.49) represents a function $f(\mathbf{z})$ as accurately as the simple Taylor series (8.19) represents a function $f(z)$, and for the same reason.

Chapter 9

Integration techniques

Equation (4.19) implies a general technique for calculating a derivative symbolically. Its counterpart (7.1), unfortunately, implies a general technique only for calculating an integral *numerically*—and even for this purpose it is imperfect; for, when it comes to adding an infinite number of infinitesimal elements, how is one actually to do the sum?

It turns out that there is no one general answer to this question. Some functions are best integrated by one technique, some by another. It is hard to guess in advance which technique might work best.

This chapter surveys several weapons of the intrepid mathematician's arsenal against the integral.

9.1 Integration by antiderivative

The simplest way to solve an integral is just to look at it, recognizing its integrand to be the derivative of something already known:¹

$$\int_a^z \frac{df}{d\tau} d\tau = f(\tau)|_a^z. \quad (9.1)$$

For instance,

$$\int_1^x \frac{1}{\tau} d\tau = \ln \tau|_1^x = \ln x.$$

One merely looks at the integrand $1/\tau$, recognizing it to be the derivative of $\ln \tau$, then directly writes down the solution $\ln \tau|_1^x$. Refer to § 7.2.

¹The notation $f(\tau)|_a^z$ or $[f(\tau)]_a^z$ means $f(z) - f(a)$.

The technique by itself is pretty limited. However, the frequent object of other integration techniques is to transform an integral into a form to which this basic technique can be applied.

Besides the essential

$$\tau^{a-1} = \frac{d}{d\tau} \left(\frac{\tau^a}{a} \right), \quad (9.2)$$

Tables 7.1, 5.2, 5.3 and 9.1 provide several further good derivatives this antiderivative technique can use.

One particular, nonobvious, useful variation on the antiderivative technique seems worth calling out specially here. If $z = \rho e^{i\phi}$, then (8.24) and (8.25) have that

$$\int_{z_1}^{z_2} \frac{dz}{z} = \ln \frac{\rho_2}{\rho_1} + i(\phi_2 - \phi_1). \quad (9.3)$$

This helps, for example, when z_1 and z_2 are real but negative numbers.

9.2 Integration by substitution

Consider the integral

$$S = \int_{x_1}^{x_2} \frac{x \, dx}{1 + x^2}.$$

This integral is not in a form one immediately recognizes. However, with the change of variable

$$u \leftarrow 1 + x^2,$$

whose differential is (by successive steps)

$$\begin{aligned} d(u) &= d(1 + x^2), \\ du &= 2x \, dx, \end{aligned}$$

the integral is

$$\begin{aligned}
 S &= \int_{x=x_1}^{x_2} \frac{x \, dx}{u} \\
 &= \int_{x=x_1}^{x_2} \frac{2x \, dx}{2u} \\
 &= \int_{u=1+x_1^2}^{1+x_2^2} \frac{du}{2u} \\
 &= \frac{1}{2} \ln u \Big|_{u=1+x_1^2}^{1+x_2^2} \\
 &= \frac{1}{2} \ln \frac{1+x_2^2}{1+x_1^2}.
 \end{aligned}$$

To check the result, we can take the derivative per § 7.5 of the final expression with respect to x_2 :

$$\begin{aligned}
 \frac{\partial}{\partial x_2} \frac{1}{2} \ln \frac{1+x_2^2}{1+x_1^2} \Big|_{x_2=x} &= \left[\frac{1}{2} \frac{\partial}{\partial x_2} \{ \ln(1+x_2^2) - \ln(1+x_1^2) \} \right]_{x_2=x} \\
 &= \frac{x}{1+x^2},
 \end{aligned}$$

which indeed has the form of the integrand we started with.

The technique is *integration by substitution*. It does not solve all integrals but it does solve many, whether alone or in combination with other techniques.

9.3 Integration by parts

Integration by parts is a curious but very broadly applicable technique which begins with the derivative product rule (4.25),

$$d(uv) = u \, dv + v \, du,$$

where $u(\tau)$ and $v(\tau)$ are functions of an independent variable τ . Reordering terms,

$$u \, dv = d(uv) - v \, du.$$

Integrating,

$$\int_{\tau=a}^b u \, dv = uv \Big|_{\tau=a}^b - \int_{\tau=a}^b v \, du. \quad (9.4)$$

Equation (9.4) is the rule of *integration by parts*.

For an example of the rule's operation, consider the integral

$$S(x) = \int_0^x \tau \cos \alpha \tau \, d\tau.$$

Unsure how to integrate this, we can begin by integrating *part* of it. We can begin by integrating the $\cos \alpha \tau \, d\tau$ part. Letting

$$\begin{aligned} u &\leftarrow \tau, \\ dv &\leftarrow \cos \alpha \tau \, d\tau, \end{aligned}$$

we find that²

$$\begin{aligned} du &= d\tau, \\ v &= \frac{\sin \alpha \tau}{\alpha}. \end{aligned}$$

According to (9.4), then,

$$S(x) = \left. \frac{\tau \sin \alpha \tau}{\alpha} \right|_0^x - \int_0^x \frac{\sin \alpha \tau}{\alpha} \, d\tau = \frac{x}{\alpha} \sin \alpha x + \cos \alpha x - 1.$$

Integration by parts is a powerful technique, but one should understand clearly what it does and does not do. It does not just integrate each part of an integral separately. It isn't that simple. What it does is to integrate one part of an integral separately—whichever part one has chosen to identify as dv —while contrarily differentiating the other part u , upon which it rewards the mathematician only with a whole new integral $\int v \, du$. The new integral may or may not be easier to integrate than was the original $\int u \, dv$. The virtue of the technique lies in that one often can find a part dv which does yield an easier $\int v \, du$. The technique is powerful for this reason.

For another kind of example of the rule's operation, consider the definite integral³

$$\Gamma(z) \equiv \int_0^\infty e^{-\tau} \tau^{z-1} \, d\tau, \quad \Re(z) > 0. \quad (9.5)$$

²The careful reader will observe that $v = (\sin \alpha \tau)/\alpha + C$ matches the chosen dv for any value of C , not just for $C = 0$. This is true. However, nothing in the integration by parts technique requires us to consider all possible v . Any convenient v suffices. In this case, we choose $v = (\sin \alpha \tau)/\alpha$.

³[28]

Letting

$$\begin{aligned} u &\leftarrow e^{-\tau}, \\ dv &\leftarrow \tau^{z-1} d\tau, \end{aligned}$$

we evidently have that

$$\begin{aligned} du &= -e^{-\tau} d\tau, \\ v &= \frac{\tau^z}{z}. \end{aligned}$$

Substituting these according to (9.4) into (9.5) yields

$$\begin{aligned} \Gamma(z) &= \left[e^{-\tau} \frac{\tau^z}{z} \right]_{\tau=0}^{\infty} - \int_0^{\infty} \left(\frac{\tau^z}{z} \right) (-e^{-\tau} d\tau) \\ &= [0 - 0] + \int_0^{\infty} \frac{\tau^z}{z} e^{-\tau} d\tau \\ &= \frac{\Gamma(z+1)}{z}. \end{aligned}$$

When written

$$\Gamma(z+1) = z\Gamma(z), \quad (9.6)$$

this is an interesting result. Since per (9.5)

$$\Gamma(1) = \int_0^{\infty} e^{-\tau} d\tau = [-e^{-\tau}]_0^{\infty} = 1,$$

it follows by induction on (9.6) that

$$(n-1)! = \Gamma(n). \quad (9.7)$$

Thus (9.5), called the *gamma function*, can be taken as an extended definition of the factorial $(z-1)!$ for all z , $\Re(z) > 0$. Integration by parts has made this finding possible.

9.4 Integration by unknown coefficients

One of the more powerful integration techniques is relatively inelegant, yet it easily cracks some integrals that give other techniques trouble. The technique is the *method of unknown coefficients*, and it is based on the antiderivative (9.1) plus intelligent guessing. It is best illustrated by example.

Consider the integral (which arises in probability theory)

$$S(x) = \int_0^x e^{-(\rho/\sigma)^2/2} \rho \, d\rho. \quad (9.8)$$

If one does not know how to solve the integral in a more elegant way, one can *guess* a likely-seeming antiderivative form, such as

$$e^{-(\rho/\sigma)^2/2} \rho = \frac{d}{d\rho} a e^{-(\rho/\sigma)^2/2},$$

where the a is an *unknown coefficient*. Having guessed, one has no guarantee that the guess is right, but see: if the guess *were* right, then the antiderivative would have the form

$$\begin{aligned} e^{-(\rho/\sigma)^2/2} \rho &= \frac{d}{d\rho} a e^{-(\rho/\sigma)^2/2} \\ &= -\frac{a\rho}{\sigma^2} e^{-(\rho/\sigma)^2/2}, \end{aligned}$$

implying that

$$a = -\sigma^2$$

(evidently the guess is right, after all). Using this value for a , one can write the specific antiderivative

$$e^{-(\rho/\sigma)^2/2} \rho = \frac{d}{d\rho} \left[-\sigma^2 e^{-(\rho/\sigma)^2/2} \right],$$

with which one can solve the integral, concluding that

$$S(x) = \left[-\sigma^2 e^{-(\rho/\sigma)^2/2} \right]_0^x = (\sigma^2) \left(1 - e^{-(x/\sigma)^2/2} \right). \quad (9.9)$$

The same technique solves differential equations, too. Consider for example the differential equation

$$dx = (Ix - P) dt, \quad x|_{t=0} = x_o, \quad x|_{t=T} = 0, \quad (9.10)$$

which conceptually represents⁴ the changing balance x of a bank loan account over time t , where I is the loan's interest rate and P is the borrower's payment rate. If it is desired to find the correct payment rate P which pays

⁴Real banks (in the author's country, at least) by law or custom actually use a needlessly more complicated formula—and not only more complicated, but mathematically slightly incorrect, too.

the loan off in the time T , then (perhaps after some bad guesses) we guess the form

$$x(t) = Ae^{\alpha t} + B,$$

where α , A and B are unknown coefficients. The guess' derivative is

$$dx = \alpha Ae^{\alpha t} dt.$$

Substituting the last two equations into (9.10) and dividing by dt yields

$$\alpha Ae^{\alpha t} = IAe^{\alpha t} + IB - P,$$

which at least is satisfied if both of the equations

$$\begin{aligned}\alpha Ae^{\alpha t} &= IAe^{\alpha t}, \\ 0 &= IB - P,\end{aligned}$$

are satisfied. Evidently good choices for α and B , then, are

$$\begin{aligned}\alpha &= I, \\ B &= \frac{P}{I}.\end{aligned}$$

Substituting these coefficients into the $x(t)$ equation above yields the general solution

$$x(t) = Ae^{It} + \frac{P}{I} \tag{9.11}$$

to (9.10). The constants A and P , we establish by applying the given *boundary conditions* $x|_{t=0} = x_o$ and $x|_{t=T} = 0$. For the former condition, (9.11) is

$$x_o = Ae^{(I)(0)} + \frac{P}{I} = A + \frac{P}{I};$$

and for the latter condition,

$$0 = Ae^{IT} + \frac{P}{I}.$$

Solving the last two equations simultaneously, we have that

$$\begin{aligned}A &= \frac{-e^{-IT}x_o}{1 - e^{-IT}}, \\ P &= \frac{Ix_o}{1 - e^{-IT}}.\end{aligned} \tag{9.12}$$

Applying these to the general solution (9.11) yields the specific solution

$$x(t) = \frac{x_o}{1 - e^{-IT}} \left(1 - e^{(I)(t-T)}\right) \quad (9.13)$$

to (9.10) meeting the boundary conditions, with the payment rate P required of the borrower given by (9.12).

The virtue of the method of unknown coefficients lies in that it permits one to try an entire family of candidate solutions at once, with the family members distinguished by the values of the coefficients. If a solution exists anywhere in the family, the method usually finds it.

The method of unknown coefficients is an elephant. Slightly inelegant the method may be, but it is pretty powerful, too—and it has surprise value (for some reason people seem not to expect it). Such are the kinds of problems the method can solve.

9.5 Integration by closed contour

We pass now from the elephant to the falcon, from the inelegant to the sublime. Consider the definite integral⁵

$$S = \int_0^\infty \frac{\tau^a}{\tau + 1} d\tau, \quad -1 < a < 0.$$

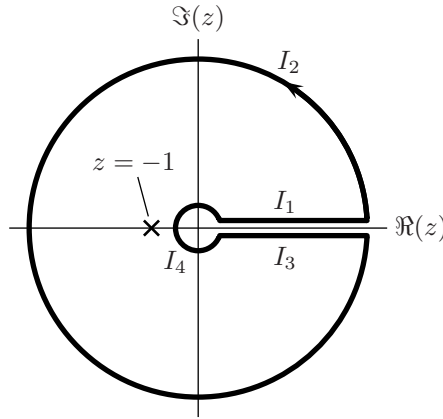
This is a hard integral. No obvious substitution, no evident factoring into parts, seems to solve the integral; but there is a way. The integrand has a pole at $\tau = -1$. Observing that τ is only a dummy integration variable, if one writes the same integral using the complex variable z in place of the real variable τ , then Cauchy's integral formula (8.29) has that integrating once counterclockwise about a closed complex contour, with the contour enclosing the pole at $z = -1$ but shutting out the branch point at $z = 0$, yields

$$I = \oint \frac{z^a}{z + 1} dz = i2\pi z^a|_{z=-1} = i2\pi \left(e^{i2\pi/2}\right)^a = i2\pi e^{i2\pi a/2}.$$

The trouble, of course, is that the integral S does not go about a closed complex contour. One can however construct a closed complex contour I of which S is a part, as in Fig 9.1. If the outer circle in the figure is of infinite

⁵[28, § 1.2]

Figure 9.1: Integration by closed contour.



radius and the inner, of infinitesimal, then the closed contour I is composed of the four parts

$$\begin{aligned} I &= I_1 + I_2 + I_3 + I_4 \\ &= (I_1 + I_3) + I_2 + I_4. \end{aligned}$$

The figure tempts one to make the mistake of writing that $I_1 = S = -I_3$, but besides being incorrect this defeats the purpose of the closed contour technique. More subtlety is needed. One must take care to interpret the four parts correctly. The integrand $z^a/(z+1)$ is multiple-valued; so, in fact, the two parts $I_1 + I_3 \neq 0$ do not cancel. The integrand has a branch point at $z = 0$, which, in passing from I_3 through I_4 to I_1 , the contour has circled. Even though z itself takes on the same values along I_3 as along I_1 , the multiple-valued integrand $z^a/(z+1)$ does not. Indeed,

$$\begin{aligned} I_1 &= \int_0^\infty \frac{(\rho e^{i0})^a}{(\rho e^{i0}) + 1} d\rho = \int_0^\infty \frac{\rho^a}{\rho + 1} d\rho = S, \\ -I_3 &= \int_0^\infty \frac{(\rho e^{i2\pi})^a}{(\rho e^{i2\pi}) + 1} d\rho = e^{i2\pi a} \int_0^\infty \frac{\rho^a}{\rho + 1} d\rho = e^{i2\pi a} S. \end{aligned}$$

Therefore,

$$\begin{aligned}
I &= I_1 + I_2 + I_3 + I_4 \\
&= (I_1 + I_3) + I_2 + I_4 \\
&= (1 - e^{i2\pi a})S + \lim_{\rho \rightarrow \infty} \int_{\phi=0}^{2\pi} \frac{z^a}{z+1} dz - \lim_{\rho \rightarrow 0} \int_{\phi=0}^{2\pi} \frac{z^a}{z+1} dz \\
&= (1 - e^{i2\pi a})S + \lim_{\rho \rightarrow \infty} \int_{\phi=0}^{2\pi} z^{a-1} dz - \lim_{\rho \rightarrow 0} \int_{\phi=0}^{2\pi} z^a dz \\
&= (1 - e^{i2\pi a})S + \lim_{\rho \rightarrow \infty} \frac{z^a}{a} \Big|_{\phi=0}^{2\pi} - \lim_{\rho \rightarrow 0} \frac{z^{a+1}}{a+1} \Big|_{\phi=0}^{2\pi}.
\end{aligned}$$

Since $a < 0$, the first limit vanishes; and because $a > -1$, the second limit vanishes, too, leaving

$$I = (1 - e^{i2\pi a})S.$$

But by Cauchy's integral formula we have already found an expression for I . Substituting this expression into the last equation yields, by successive steps,

$$\begin{aligned}
i2\pi e^{i2\pi a/2} &= (1 - e^{i2\pi a})S, \\
S &= \frac{i2\pi e^{i2\pi a/2}}{1 - e^{i2\pi a}}, \\
S &= \frac{i2\pi}{e^{-i2\pi a/2} - e^{i2\pi a/2}}, \\
S &= -\frac{2\pi/2}{\sin(2\pi a/2)}.
\end{aligned}$$

That is,

$$\int_0^\infty \frac{\tau^a}{\tau+1} d\tau = -\frac{2\pi/2}{\sin(2\pi a/2)}, \quad -1 < a < 0, \quad (9.14)$$

an astonishing result.⁶

Another example⁷ is

$$T = \int_0^{2\pi} \frac{d\theta}{1 + a \cos \theta}, \quad \Im(a) = 0, \quad |\Re(a)| < 1.$$

⁶So astonishing is the result, that one is unlikely to believe it at first encounter. However, straightforward (though computationally highly inefficient) numerical integration per (7.1) confirms the result, as the interested reader and his computer can check. Such results vindicate the effort we have spent in deriving Cauchy's integral formula (8.29).

⁷[25]

As in the previous example, here again the contour is not closed. The previous example closed the contour by extending it, excluding the branch point. In this example there is no branch point to exclude, nor need one extend the contour. Rather, one changes the variable

$$z \leftarrow e^{i\theta}$$

and takes advantage of the fact that z , unlike θ , *begins and ends the integration at the same point*. One thus obtains the equivalent integral

$$\begin{aligned} T &= \oint \frac{dz/iz}{1 + (a/2)(z + 1/z)} = -\frac{i2}{a} \oint \frac{dz}{z^2 + 2z/a + 1} \\ &= -\frac{i2}{a} \oint \frac{dz}{\left[z - \left(-1 + \sqrt{1 - a^2}\right)/a\right] \left[z - \left(-1 - \sqrt{1 - a^2}\right)/a\right]}, \end{aligned}$$

whose contour is the unit circle in the Argand plane. The integrand evidently has poles at

$$z = \frac{-1 \pm \sqrt{1 - a^2}}{a},$$

whose magnitudes are such that

$$|z|^2 = \frac{2 - a^2 \mp 2\sqrt{1 - a^2}}{a^2}.$$

One of the two magnitudes is less than unity and one is greater, meaning that one of the two poles lies within the contour and one lies without, as is

seen by the successive steps⁸

$$\begin{aligned}
 a^2 &< 1, \\
 0 &< 1 - a^2, \\
 (-a^2)(0) &> (-a^2)(1 - a^2), \\
 0 &> -a^2 + a^4, \\
 1 - a^2 &> 1 - 2a^2 + a^4, \\
 1 - a^2 &> (1 - a^2)^2, \\
 \sqrt{1 - a^2} &> 1 - a^2, \\
 -\sqrt{1 - a^2} &< -(1 - a^2) < \sqrt{1 - a^2}, \\
 1 - \sqrt{1 - a^2} &< a^2 < 1 + \sqrt{1 - a^2}, \\
 2 - 2\sqrt{1 - a^2} &< 2a^2 < 2 + 2\sqrt{1 - a^2}, \\
 2 - a^2 - 2\sqrt{1 - a^2} &< a^2 < 2 - a^2 + 2\sqrt{1 - a^2}, \\
 \frac{2 - a^2 - 2\sqrt{1 - a^2}}{a^2} &< 1 < \frac{2 - a^2 + 2\sqrt{1 - a^2}}{a^2}.
 \end{aligned}$$

Per Cauchy's integral formula (8.29), integrating about the pole within the contour yields

$$T = i2\pi \frac{-i2/a}{z - \left(-1 - \sqrt{1 - a^2}\right)/a} \Big|_{z=(-1+\sqrt{1-a^2})/a} = \frac{2\pi}{\sqrt{1 - a^2}}.$$

Observe that by means of a complex variable of integration, each example has indirectly evaluated an integral whose integrand is purely real. If it seems unreasonable to the reader to expect so flamboyant a technique actually to work, this seems equally unreasonable to the writer—but work it does, nevertheless. It is a great technique.

The technique, *integration by closed contour*, is found in practice to solve many integrals other techniques find almost impossible to crack. The key to making the technique work lies in closing a contour one knows how to treat. The robustness of the technique lies in that any contour of any shape will work, so long as the contour encloses appropriate poles in the Argand domain plane while shutting branch points out.

⁸These steps are perhaps best read from bottom to top. See Ch. 6's footnote 14.

9.6 Integration by partial-fraction expansion

This section treats integration by partial-fraction expansion. It introduces the expansion itself first.⁹

9.6.1 Partial-fraction expansion

Consider the function

$$f(z) = \frac{-4}{z-1} + \frac{5}{z-2}.$$

Combining the two fractions over a common denominator¹⁰ yields

$$f(z) = \frac{z+3}{(z-1)(z-2)}.$$

Of the two forms, the former is probably the more amenable to analysis. For example, using (9.3),

$$\begin{aligned} \int_{-1}^0 f(\tau) d\tau &= \int_{-1}^0 \frac{-4}{\tau-1} d\tau + \int_{-1}^0 \frac{5}{\tau-2} d\tau \\ &= [-4 \ln(1-\tau) + 5 \ln(2-\tau)]_{-1}^0. \end{aligned}$$

The trouble is that one is not always given the function in the amenable form.

Given a *rational function*

$$f(z) = \frac{\sum_{k=0}^{N-1} b_k z^k}{\prod_{j=1}^N (z - \alpha_j)} \quad (9.15)$$

in which no two of the several poles α_j are the same, the *partial-fraction expansion* has the form

$$f(z) = \sum_{k=1}^N \frac{A_k}{z - \alpha_k}, \quad (9.16)$$

where multiplying each fraction of (9.16) by

$$\frac{\left[\prod_{j=1}^N (z - \alpha_j) \right] / (z - \alpha_k)}{\left[\prod_{j=1}^N (z - \alpha_j) \right] / (z - \alpha_k)}$$

⁹[32, Appendix F][21, §§ 2.7 and 10.12]

¹⁰Terminology (you probably knew this already): A *fraction* is the ratio of two numbers or expressions B/A . In the fraction, B is the *numerator* and A is the *denominator*. The *quotient* is $Q = B/A$.

puts the several fractions over a common denominator, yielding (9.15). Dividing (9.15) by (9.16) gives the ratio

$$1 = \frac{\sum_{k=0}^{N-1} b_k z^k}{\prod_{j=1}^N (z - \alpha_j)} \bigg/ \sum_{k=1}^N \frac{A_k}{z - \alpha_k}.$$

In the immediate neighborhood of $z = \alpha_m$, the m th term $A_m/(z - \alpha_m)$ dominates the summation of (9.16). Hence,

$$1 = \lim_{z \rightarrow \alpha_m} \frac{\sum_{k=0}^{N-1} b_k z^k}{\prod_{j=1}^N (z - \alpha_j)} \bigg/ \frac{A_m}{z - \alpha_m}.$$

Rearranging factors, we have that

$$A_m = \left. \frac{\sum_{k=0}^{N-1} b_k z^k}{\left[\prod_{j=1}^N (z - \alpha_j) \right] / (z - \alpha_m)} \right|_{z=\alpha_m} = \lim_{z \rightarrow \alpha_m} [(z - \alpha_m)f(z)], \quad (9.17)$$

where A_m , the value of $f(z)$ with the pole canceled, is called the *residue* of $f(z)$ at the pole $z = \alpha_m$. Equations (9.16) and (9.17) together give the partial-fraction expansion of (9.15)'s rational function $f(z)$.

9.6.2 Repeated poles

The weakness of the partial-fraction expansion of § 9.6.1 is that it cannot directly handle repeated poles. That is, if $\alpha_i = \alpha_j$, $i \neq j$, then the residue formula (9.17) finds an uncanceled pole remaining in its denominator and thus fails for $A_i = A_j$ (it still works for the other A_m). The conventional way to expand a fraction with repeated poles is presented in § 9.6.5 below; but because at least to this writer that way does not lend much applied insight, the present subsection treats the matter in a different way. Here, we *separate the poles*.

Consider the function

$$g(z) = \sum_{k=0}^{N-1} \frac{C e^{i2\pi k/N}}{z - \epsilon e^{i2\pi k/N}}, \quad N > 1, \quad 0 < \epsilon \ll 1, \quad (9.18)$$

where C is a real-valued constant. This function evidently has a small circle of poles in the Argand plane at $\alpha_k = \epsilon e^{i2\pi k/N}$. Factoring,

$$g(z) = \frac{C}{z} \sum_{k=0}^{N-1} \frac{e^{i2\pi k/N}}{1 - \epsilon e^{i2\pi k/N}/z}.$$

Using (2.34) to expand the fraction,

$$\begin{aligned}
 g(z) &= \frac{C}{z} \sum_{k=0}^{N-1} \left[e^{i2\pi k/N} \sum_{j=0}^{\infty} \left(\frac{\epsilon e^{i2\pi k/N}}{z} \right)^j \right] \\
 &= C \sum_{k=0}^{N-1} \sum_{j=1}^{\infty} \frac{\epsilon^{j-1} e^{i2\pi jk/N}}{z^j} \\
 &= C \sum_{j=1}^{\infty} \frac{\epsilon^{j-1}}{z^j} \sum_{k=0}^{N-1} \left(e^{i2\pi j/N} \right)^k.
 \end{aligned}$$

But¹¹

$$\sum_{k=0}^{N-1} \left(e^{i2\pi j/N} \right)^k = \begin{cases} N & \text{if } j = mN, \\ 0 & \text{otherwise,} \end{cases}$$

so

$$g(z) = NC \sum_{m=1}^{\infty} \frac{\epsilon^{mN-1}}{z^{mN}}.$$

For $|z| \gg \epsilon$ —that is, except in the immediate neighborhood of the small circle of poles—the first term of the summation dominates. Hence,

$$g(z) \approx NC \frac{\epsilon^{N-1}}{z^N}, \quad |z| \gg \epsilon.$$

Having achieved this approximation, if we strategically choose

$$C = \frac{1}{N\epsilon^{N-1}},$$

then

$$g(z) \approx \frac{1}{z^N}, \quad |z| \gg \epsilon.$$

But given the chosen value of C , (9.18) is

$$g(z) = \frac{1}{N\epsilon^{N-1}} \sum_{k=0}^{N-1} \frac{e^{i2\pi k/N}}{z - \epsilon e^{i2\pi k/N}}, \quad N > 1, \quad 0 < \epsilon \ll 1.$$

¹¹If you don't see why, then for $N = 8$ and $j = 3$ plot the several $(e^{i2\pi j/N})^k$ in the Argand plane. Do the same for $j = 2$ then $j = 8$. Only in the $j = 8$ case do the terms add coherently; in the other cases they cancel.

This effect—reinforcing when $j = nN$, canceling otherwise—is a classic manifestation of *Parseval's principle*.

Joining the last two equations together, changing $z - z_o \leftarrow z$, and writing more formally, we have that

$$\frac{1}{(z - z_o)^N} = \lim_{\epsilon \rightarrow 0} \frac{1}{N\epsilon^{N-1}} \sum_{k=0}^{N-1} \frac{e^{i2\pi k/N}}{z - [z_o + \epsilon e^{i2\pi k/N}]}, \quad N > 1. \quad (9.19)$$

The significance of (9.19) is that it lets one replace an N -fold pole with a small circle of ordinary poles, which per § 9.6.1 we already know how to handle. Notice incidentally that $1/N\epsilon^{N-1}$ is a large number not a small. The poles are close together but very strong.

An example to illustrate the technique, separating a double pole:

$$\begin{aligned} f(z) &= \frac{z^2 - z + 6}{(z - 1)^2(z + 2)} \\ &= \lim_{\epsilon \rightarrow 0} \frac{z^2 - z + 6}{(z - [1 + \epsilon e^{i2\pi(0)/2}])(z - [1 + \epsilon e^{i2\pi(1)/2}])(z + 2)} \\ &= \lim_{\epsilon \rightarrow 0} \frac{z^2 - z + 6}{(z - [1 + \epsilon])(z - [1 - \epsilon])(z + 2)} \\ &= \lim_{\epsilon \rightarrow 0} \left\{ \left(\frac{1}{z - [1 + \epsilon]} \right) \left[\frac{z^2 - z + 6}{(z - [1 - \epsilon])(z + 2)} \right]_{z=1+\epsilon} \right. \\ &\quad \left. + \left(\frac{1}{z - [1 - \epsilon]} \right) \left[\frac{z^2 - z + 6}{(z - [1 + \epsilon])(z + 2)} \right]_{z=1-\epsilon} \right. \\ &\quad \left. + \left(\frac{1}{z + 2} \right) \left[\frac{z^2 - z + 6}{(z - [1 + \epsilon])(z - [1 - \epsilon])} \right]_{z=-2} \right\} \\ &= \lim_{\epsilon \rightarrow 0} \left\{ \left(\frac{1}{z - [1 + \epsilon]} \right) \left[\frac{6 + \epsilon}{6\epsilon + 2\epsilon^2} \right] \right. \\ &\quad \left. + \left(\frac{1}{z - [1 - \epsilon]} \right) \left[\frac{6 - \epsilon}{-6\epsilon + 2\epsilon^2} \right] \right. \\ &\quad \left. + \left(\frac{1}{z + 2} \right) \left[\frac{0xC}{9} \right] \right\} \\ &= \lim_{\epsilon \rightarrow 0} \left\{ \frac{1/\epsilon - 1/6}{z - [1 + \epsilon]} + \frac{-1/\epsilon - 1/6}{z - [1 - \epsilon]} + \frac{4/3}{z + 2} \right\} \\ &= \lim_{\epsilon \rightarrow 0} \left\{ \frac{1/\epsilon}{z - [1 + \epsilon]} + \frac{-1/\epsilon}{z - [1 - \epsilon]} + \frac{-1/3}{z - 1} + \frac{4/3}{z + 2} \right\}. \end{aligned}$$

Notice how the calculation has discovered an additional, single pole at $z = 1$, the pole hiding under dominant, double pole there.

9.6.3 Integrating a rational function

If one can find the poles of a rational function of the form (9.15), then one can use (9.16) and (9.17)—and, if needed, (9.19)—to expand the function into a sum of partial fractions, each of which one can integrate individually.

Continuing the example of § 9.6.2, for $0 \leq x < 1$,

$$\begin{aligned}
 \int_0^x f(\tau) d\tau &= \int_0^x \frac{\tau^2 - \tau + 6}{(\tau - 1)^2(\tau + 2)} d\tau \\
 &= \lim_{\epsilon \rightarrow 0} \int_0^x \left\{ \frac{1/\epsilon}{\tau - [1 + \epsilon]} + \frac{-1/\epsilon}{\tau - [1 - \epsilon]} + \frac{-1/3}{\tau - 1} + \frac{4/3}{\tau + 2} \right\} d\tau \\
 &= \lim_{\epsilon \rightarrow 0} \left\{ \frac{1}{\epsilon} \ln([1 + \epsilon] - \tau) - \frac{1}{\epsilon} \ln([1 - \epsilon] - \tau) \right. \\
 &\quad \left. - \frac{1}{3} \ln(1 - \tau) + \frac{4}{3} \ln(\tau + 2) \right\}_0^x \\
 &= \lim_{\epsilon \rightarrow 0} \left\{ \frac{1}{\epsilon} \ln \left(\frac{[1 + \epsilon] - \tau}{[1 - \epsilon] - \tau} \right) - \frac{1}{3} \ln(1 - \tau) + \frac{4}{3} \ln(\tau + 2) \right\}_0^x \\
 &= \lim_{\epsilon \rightarrow 0} \left\{ \frac{1}{\epsilon} \ln \left(\frac{[1 - \tau] + \epsilon}{[1 - \tau] - \epsilon} \right) - \frac{1}{3} \ln(1 - \tau) + \frac{4}{3} \ln(\tau + 2) \right\}_0^x \\
 &= \lim_{\epsilon \rightarrow 0} \left\{ \frac{1}{\epsilon} \ln \left(1 + \frac{2\epsilon}{1 - \tau} \right) - \frac{1}{3} \ln(1 - \tau) + \frac{4}{3} \ln(\tau + 2) \right\}_0^x \\
 &= \lim_{\epsilon \rightarrow 0} \left\{ \frac{1}{\epsilon} \left(\frac{2\epsilon}{1 - \tau} \right) - \frac{1}{3} \ln(1 - \tau) + \frac{4}{3} \ln(\tau + 2) \right\}_0^x \\
 &= \lim_{\epsilon \rightarrow 0} \left\{ \frac{2}{1 - \tau} - \frac{1}{3} \ln(1 - \tau) + \frac{4}{3} \ln(\tau + 2) \right\}_0^x \\
 &= \frac{2}{1 - x} - 2 - \frac{1}{3} \ln(1 - x) + \frac{4}{3} \ln \left(\frac{x + 2}{2} \right).
 \end{aligned}$$

To check (§ 7.5) that the result is correct, we can take the derivative of the final expression:

$$\begin{aligned}
 &\left[\frac{d}{dx} \left\{ \frac{2}{1 - x} - 2 - \frac{1}{3} \ln(1 - x) + \frac{4}{3} \ln \left(\frac{x + 2}{2} \right) \right\} \right]_{x=\tau} \\
 &= \frac{2}{(\tau - 1)^2} + \frac{-1/3}{\tau - 1} + \frac{4/3}{\tau + 2} \\
 &= \frac{\tau^2 - \tau + 6}{(\tau - 1)^2(\tau + 2)},
 \end{aligned}$$

which indeed has the form of the integrand we started with, confirming the

result. (Notice incidentally how much easier it is symbolically to differentiate than to integrate!)

9.6.4 The derivatives of a rational function

Not only the integral of a rational function interests us; its derivatives interest us, too. One needs no special technique to compute such derivatives, of course, but the derivatives do bring some noteworthy properties.

First of interest is the property that a function in the general rational form

$$\Phi(w) = \frac{w^p h_0(w)}{g(w)}, \quad g(0) \neq 0, \quad (9.20)$$

enjoys derivatives in the general rational form

$$\frac{d^k \Phi}{dw^k} = \frac{w^{p-k} h_k(w)}{[g(w)]^{k+1}}, \quad 0 \leq k \leq p, \quad (9.21)$$

where g and h_k are polynomials in nonnegative powers of w . The property is proved by induction. When $k = 0$, (9.21) is (9.20), so (9.21) is good at least for this case. Then, if (9.21) holds for $k = i - 1$,

$$\begin{aligned} \frac{d^i \Phi}{dw^i} &= \frac{d}{dw} \left[\frac{d^{i-1} \Phi}{dw^{i-1}} \right] = \frac{d}{dw} \left[\frac{w^{p-i+1} h_{i-1}(w)}{[g(w)]^i} \right] = \frac{w^{p-i} h_i(w)}{[g(w)]^{i+1}}, \\ h_i(w) &\equiv w g \frac{dh_{i-1}}{dw} - i w h_{i-1} \frac{dg}{dw} + (p - i + 1) g h_{i-1}, \quad 0 < i \leq p, \end{aligned}$$

which makes h_i (like h_{i-1}) a polynomial in nonnegative powers of w . By induction on this basis, (9.21) holds for all $0 \leq k \leq p$, as was to be demonstrated.

A related property is that

$$\left. \frac{d^k \Phi}{dw^k} \right|_{w=0} = 0 \quad \text{for } 0 \leq k < p. \quad (9.22)$$

That is, the function and its first $p - 1$ derivatives are all zero at $w = 0$. The reason is that (9.21)'s denominator is $[g(w)]^{k+1} \neq 0$, whereas its numerator has a $w^{p-k} = 0$ factor, when $0 \leq k < p$ and $w = 0$.

9.6.5 Repeated poles (the conventional technique)

Though the technique of §§ 9.6.2 and 9.6.3 affords extra insight, it is not the conventional technique to expand in partial fractions a rational function

having a repeated pole. The conventional technique is worth learning not only because it is conventional but also because it is usually quicker to apply in practice. This subsection derives it.

A rational function with repeated poles,

$$\begin{aligned} f(z) &= \frac{\sum_{k=0}^{N-1} b_k z^k}{\prod_{j=1}^M (z - \alpha_j)^{p_j}}, \\ N &\equiv \sum_{j=1}^M p_j, \\ p_j &\geq 0, \\ \alpha_{j'} &\neq \alpha_j \text{ if } j' \neq j, \end{aligned} \tag{9.23}$$

where j, k, M, N and the several p_j are integers, cannot be expanded solely in the first-order fractions of § 9.6.1, but can indeed be expanded if higher-order fractions are allowed:

$$f(z) = \sum_{j=1}^M \sum_{\ell=0}^{p_j-1} \frac{A_{j\ell}}{(z - \alpha_j)^{p_j-\ell}}. \tag{9.24}$$

What the partial-fraction expansion (9.24) lacks are the values of its several coefficients $A_{j\ell}$.

One can determine the coefficients with respect to one (possibly repeated) pole at a time. To determine them with respect to the p_m -fold pole at $z = \alpha_m$, $1 \leq m \leq M$, one multiplies (9.24) by $(z - \alpha_m)^{p_m}$ to obtain the form

$$(z - \alpha_m)^{p_m} f(z) = \sum_{\substack{j=1, \\ j \neq m}}^M \sum_{\ell=0}^{p_j-1} \frac{(A_{j\ell})(z - \alpha_m)^{p_m}}{(z - \alpha_j)^{p_j-\ell}} + \sum_{\ell=0}^{p_m-1} (A_{m\ell})(z - \alpha_m)^\ell.$$

But (9.22) with $w = z - \alpha_m$ reveals the double summation and its first $p_m - 1$ derivatives all to be null at $z = \alpha_m$; that is,

$$\left. \frac{d^k}{dz^k} \sum_{\substack{j=1, \\ j \neq m}}^M \sum_{\ell=0}^{p_j-1} \frac{(A_{j\ell})(z - \alpha_m)^{p_m}}{(z - \alpha_j)^{p_j-\ell}} \right|_{z=\alpha_m} = 0, \quad 0 \leq k < p_m;$$

so, the $(z - \alpha_m)^{p_m} f(z)$ equation's k th derivative reduces at that point to

$$\begin{aligned} \left. \frac{d^k}{dz^k} [(z - \alpha_m)^{p_m} f(z)] \right|_{z=\alpha_m} &= \sum_{\ell=0}^{p_m-1} \frac{d^k}{dz^k} [(A_{m\ell})(z - \alpha_m)^\ell] \Big|_{z=\alpha_m} \\ &= k! A_{mk}, \quad 0 \leq k < p_m. \end{aligned}$$

Changing $j \leftarrow m$ and $\ell \leftarrow k$ and solving for $A_{j\ell}$ then produces the coefficients

$$A_{j\ell} = \left(\frac{1}{\ell!} \right) \frac{d^\ell}{dz^\ell} [(z - \alpha_j)^{p_j} f(z)] \Big|_{z=\alpha_j}, \quad 0 \leq \ell < p_j, \quad (9.25)$$

to weight the expansion (9.24)'s partial fractions. In case of a repeated pole, these coefficients evidently depend not only on the residual function itself but also on its several derivatives, one derivative per repetition of the pole.

9.6.6 The existence and uniqueness of solutions

Equation (9.25) has solved (9.23) and (9.24). A professional mathematician might object however that it has done so without first proving that a unique solution actually exists.

Comes from us the reply, "Why should we prove that a solution exists, once we have actually found it?"

Ah, but the professional's point is that we have found the solution only if in fact it does exist, and uniquely; otherwise what we have *found* is a phantom. A careful review of § 9.6.5's logic discovers no guarantee that all of (9.25)'s coefficients actually come from the same expansion. Maybe there exist two distinct expansions, and some of the coefficients come from the one, some from the other. On the other hand, maybe there exists no expansion at all, in which event it is not even clear what (9.25) means.

"But these are quibbles, cavils and nitpicks!" we are inclined to grumble. "The present book is a book of applied mathematics."

Well, yes, but on this occasion let us nonetheless follow the professional's line of reasoning, if only a short way.

Uniqueness is proved by positing two solutions

$$f(z) = \sum_{j=1}^M \sum_{\ell=0}^{p_j-1} \frac{A_{j\ell}}{(z - \alpha_j)^{p_j-\ell}} = \sum_{j=1}^M \sum_{\ell=0}^{p_j-1} \frac{B_{j\ell}}{(z - \alpha_j)^{p_j-\ell}}$$

and computing the difference

$$\sum_{j=1}^M \sum_{\ell=0}^{p_j-1} \frac{B_{j\ell} - A_{j\ell}}{(z - \alpha_j)^{p_j-\ell}}$$

between them. Logically this difference must be zero for all z if the two solutions are actually to represent the same function $f(z)$. This however is seen to be possible only if $B_{j\ell} = A_{j\ell}$ for each (j, ℓ) . Therefore, the two solutions are one and the same.

Existence comes of combining the several fractions of (9.24) over a common denominator and comparing the resulting numerator against the numerator of (9.23). Each coefficient b_k is seen thereby to be a linear combination of the several $A_{j\ell}$, where the combination's weights depend solely on the locations α_j and multiplicities p_j of $f(z)$'s several poles. From the N coefficients b_k and the N coefficients $A_{j\ell}$, an $N \times N$ system of N linear equations in N unknowns results—which might for example (if, say, $N = 3$) look like

$$\begin{aligned} b_0 &= -2A_{00} + A_{01} + 3A_{10}, \\ b_1 &= A_{00} + A_{01} + A_{10}, \\ b_2 &= 2A_{01} - 5A_{10}. \end{aligned}$$

We shall show in Chs. 11, 12 and 14 that when such a system has no solution, there always exist an alternate set of b_k for which the same system has multiple solutions. But uniqueness, which we have already established, forbids such multiple solutions in all cases. Therefore it is not possible for the system to have no solution—which is to say, the solution necessarily exists.

We shall not often in this book prove existence and uniqueness explicitly, but such proofs when desired tend to fit the pattern outlined here.

9.7 Frullani's integral

One occasionally meets an integral of the form

$$S = \int_0^\infty \frac{f(b\tau) - f(a\tau)}{\tau} d\tau,$$

where a and b are real, positive coefficients and $f(\tau)$ is an arbitrary complex expression in τ . Such an integral, one wants to split in two as $\int [f(b\tau)/\tau] d\tau - \int [f(a\tau)/\tau] d\tau$; but if $f(0^+) \neq f(+\infty)$, one cannot, because each half-integral

alone diverges. Nonetheless, splitting the integral in two is the right idea, provided that one first relaxes the limits of integration as

$$S = \lim_{\epsilon \rightarrow 0^+} \left\{ \int_{\epsilon}^{1/\epsilon} \frac{f(b\tau)}{\tau} d\tau - \int_{\epsilon}^{1/\epsilon} \frac{f(a\tau)}{\tau} d\tau \right\}.$$

Changing $\sigma \leftarrow b\tau$ in the left integral and $\sigma \leftarrow a\tau$ in the right yields

$$\begin{aligned} S &= \lim_{\epsilon \rightarrow 0^+} \left\{ \int_{b\epsilon}^{b/\epsilon} \frac{f(\sigma)}{\sigma} d\sigma - \int_{a\epsilon}^{a/\epsilon} \frac{f(\sigma)}{\sigma} d\sigma \right\} \\ &= \lim_{\epsilon \rightarrow 0^+} \left\{ \int_{a\epsilon}^{b\epsilon} \frac{-f(\sigma)}{\sigma} d\sigma + \int_{b\epsilon}^{a/\epsilon} \frac{f(\sigma) - f(\sigma)}{\sigma} d\sigma + \int_{a/\epsilon}^{b/\epsilon} \frac{f(\sigma)}{\sigma} d\sigma \right\} \\ &= \lim_{\epsilon \rightarrow 0^+} \left\{ \int_{a/\epsilon}^{b/\epsilon} \frac{f(\sigma)}{\sigma} d\sigma - \int_{a\epsilon}^{b\epsilon} \frac{f(\sigma)}{\sigma} d\sigma \right\} \end{aligned}$$

(here on the face of it, we have split the integration as though $a \leq b$, but in fact it does not matter which of a and b is the greater, as is easy to verify). So long as each of $f(\epsilon)$ and $f(1/\epsilon)$ approaches a constant value as ϵ vanishes, this is

$$\begin{aligned} S &= \lim_{\epsilon \rightarrow 0^+} \left\{ f(+\infty) \int_{a/\epsilon}^{b/\epsilon} \frac{d\sigma}{\sigma} - f(0^+) \int_{a\epsilon}^{b\epsilon} \frac{d\sigma}{\sigma} \right\} \\ &= \lim_{\epsilon \rightarrow 0^+} \left\{ f(+\infty) \ln \frac{b/\epsilon}{a/\epsilon} - f(0^+) \ln \frac{b\epsilon}{a\epsilon} \right\} \\ &= [f(\tau)]_0^\infty \ln \frac{b}{a}. \end{aligned}$$

Thus we have *Frullani's integral*,

$$\int_0^\infty \frac{f(b\tau) - f(a\tau)}{\tau} d\tau = [f(\tau)]_0^\infty \ln \frac{b}{a}, \quad (9.26)$$

which, if a and b are both real and positive, works for any $f(\tau)$ which has definite $f(0^+)$ and $f(+\infty)$.¹²

¹²[28, § 1.3][1, § 2.6.1][45, “Frullani's integral”]

9.8 Integrating products of exponentials, powers and logarithms

The products $\exp(\alpha\tau)\tau^n$ and $\tau^{a-1}\ln\tau$ tend to arise¹³ among other places in integrands related to special functions ([chapter not yet written]). The two occur often enough to merit investigation here.

Concerning $\exp(\alpha\tau)\tau^n$, by § 9.4's method of unknown coefficients we guess its antiderivative to fit the form

$$\begin{aligned}\exp(\alpha\tau)\tau^n &= \frac{d}{d\tau} \sum_{k=0}^n a_k \exp(\alpha\tau)\tau^k \\ &= \sum_{k=0}^n \alpha a_k \exp(\alpha\tau)\tau^k + \sum_{k=1}^n \frac{a_k}{k} \exp(\alpha\tau)\tau^{k-1} \\ &= \alpha a_n \exp(\alpha\tau)\tau^n + \sum_{k=0}^{n-1} \left(\alpha a_k + \frac{a_{k+1}}{k+1} \right) \exp(\alpha\tau)\tau^k.\end{aligned}$$

If so, then evidently

$$\begin{aligned}a_n &= \frac{1}{\alpha}; \\ a_k &= -\frac{a_{k+1}}{(k+1)(\alpha)}, \quad 0 \leq k < n.\end{aligned}$$

That is,

$$a_k = \frac{1}{\alpha \prod_{j=k+1}^n (-j\alpha)} = \frac{(-)^{n-k}}{(n!/k!) \alpha^{n-k+1}}, \quad 0 \leq k \leq n.$$

Therefore,¹⁴

$$\exp(\alpha\tau)\tau^n = \frac{d}{d\tau} \sum_{k=0}^n \frac{(-)^{n-k}}{(n!/k!) \alpha^{n-k+1}} \exp(\alpha\tau)\tau^k, \quad n \geq 0, \alpha \neq 0. \quad (9.27)$$

The right form to guess for the antiderivative of $\tau^{a-1}\ln\tau$ is less obvious. Remembering however § 5.3's observation that $\ln\tau$ is of zeroth order in τ , after maybe some false tries we eventually do strike the right form

$$\begin{aligned}\tau^{a-1}\ln\tau &= \frac{d}{d\tau} \tau^a [B \ln\tau + C] \\ &= \tau^{a-1} [aB \ln\tau + (B + aC)],\end{aligned}$$

¹³One could write the latter product more generally as $\tau^{a-1}\ln\beta\tau$. According to Table 2.5, however, $\ln\beta\tau = \ln\beta + \ln\tau$; wherein $\ln\beta$ is just a constant.

¹⁴[37, Appendix 2, eqn. 73]

Table 9.1: Antiderivatives of products of exponentials, powers and logarithms.

$$\begin{aligned}\exp(\alpha\tau)\tau^n &= \frac{d}{d\tau} \sum_{k=0}^n \frac{(-)^{n-k}}{(n!/k!)\alpha^{n-k+1}} \exp(\alpha\tau)\tau^k, \quad n \geq 0, \alpha \neq 0 \\ \tau^{a-1} \ln \tau &= \frac{d}{d\tau} \frac{\tau^a}{a} \left(\ln \tau - \frac{1}{a} \right), \quad a \neq 0 \\ \frac{\ln \tau}{\tau} &= \frac{d}{d\tau} \frac{(\ln \tau)^2}{2}\end{aligned}$$

which demands that $B = 1/a$ and that $C = -1/a^2$. Therefore,¹⁵

$$\tau^{a-1} \ln \tau = \frac{d}{d\tau} \frac{\tau^a}{a} \left(\ln \tau - \frac{1}{a} \right), \quad a \neq 0. \quad (9.28)$$

Antiderivatives of terms like $\tau^{a-1}(\ln \tau)^2$, $\exp(\alpha\tau)\tau^n \ln \tau$ and so on can be computed in like manner as the need arises.

Equation (9.28) fails when $a = 0$, but in this case with a little imagination the antiderivative is not hard to guess:

$$\frac{\ln \tau}{\tau} = \frac{d}{d\tau} \frac{(\ln \tau)^2}{2}. \quad (9.29)$$

If (9.29) seemed hard to guess nevertheless, then l'Hôpital's rule (4.29), applied to (9.28) as $a \rightarrow 0$, with the observation from (2.41) that

$$\tau^a = \exp(a \ln \tau), \quad (9.30)$$

would yield the same (9.29).

Table 9.1 summarizes.

9.9 Integration by Taylor series

With sufficient cleverness the techniques of the foregoing sections solve many, many integrals. But not all. When all else fails, as sometimes it does, the Taylor series of Ch. 8 and the antiderivative of § 9.1 together offer a concise,

¹⁵[37, Appendix 2, eqn. 74]

practical way to integrate some functions, at the price of losing the functions' known closed analytic forms. For example,

$$\begin{aligned}
 \int_0^x \exp\left(-\frac{\tau^2}{2}\right) d\tau &= \int_0^x \sum_{k=0}^{\infty} \frac{(-\tau^2/2)^k}{k!} d\tau \\
 &= \int_0^x \sum_{k=0}^{\infty} \frac{(-)^k \tau^{2k}}{2^k k!} d\tau \\
 &= \left[\sum_{k=0}^{\infty} \frac{(-)^k \tau^{2k+1}}{(2k+1)2^k k!} \right]_0^x \\
 &= \sum_{k=0}^{\infty} \frac{(-)^k x^{2k+1}}{(2k+1)2^k k!} = (x) \sum_{k=0}^{\infty} \frac{1}{2k+1} \prod_{j=1}^k \frac{-x^2}{2j}.
 \end{aligned}$$

The result is no function one recognizes; it is just a series. This is not necessarily bad, however. After all, when a Taylor series from Table 8.1 is used to calculate $\sin z$, then $\sin z$ is just a series, too. The series above converges just as accurately and just as fast.

Sometimes it helps to give the series a name like

$$\text{myf } z \equiv \sum_{k=0}^{\infty} \frac{(-)^k z^{2k+1}}{(2k+1)2^k k!} = (z) \sum_{k=0}^{\infty} \frac{1}{2k+1} \prod_{j=1}^k \frac{-z^2}{2j}.$$

Then,

$$\int_0^x \exp\left(-\frac{\tau^2}{2}\right) d\tau = \text{myf } x.$$

The $\text{myf } z$ is no less a function than $\sin z$ is; it's just a function you hadn't heard of before. You can plot the function, or take its derivative

$$\frac{d}{d\tau} \text{myf } \tau = \exp\left(-\frac{\tau^2}{2}\right),$$

or calculate its value, or do with it whatever else one does with functions. It works just the same.

Chapter 10

Cubics and quartics

Under the heat of noonday, between the hard work of the morning and the heavy lifting of the afternoon, one likes to lay down one's burden and rest a spell in the shade. Chapters 2 through 9 have established the applied mathematical foundations upon which coming chapters will build; and Ch. 11, hefting the weighty topic of the matrix, will indeed begin to build on those foundations. But in this short chapter which rests between, we shall refresh ourselves with an interesting but lighter mathematical topic: the topic of cubics and quartics.

The expression

$$z + a_0$$

is a *linear* polynomial, the lone root $z = -a_0$ of which is plain to see. The *quadratic* polynomial

$$z^2 + a_1z + a_0$$

has of course two roots, which though not plain to see the quadratic formula (2.2) extracts with little effort. So much algebra has been known since antiquity. The roots of higher-order polynomials, the Newton-Raphson iteration (4.30) locates swiftly, but that is an approximate iteration rather than an exact formula like (2.2), and as we have seen in § 4.8 it can occasionally fail to converge. One would prefer an actual formula to extract the roots.

No general formula to extract the roots of the n th-order polynomial seems to be known.¹ However, to extract the roots of the *cubic* and *quartic* polynomials

$$\begin{aligned} z^3 + a_2z^2 + a_1z + a_0, \\ z^4 + a_3z^3 + a_2z^2 + a_1z + a_0, \end{aligned}$$

¹Refer to Ch. 6's footnote 8.

though the ancients never discovered how, formulas do exist. The 16th-century algebraists Ferrari, Vieta, Tartaglia and Cardano have given us the clever technique. This chapter explains.²

10.1 Vieta's transform

There is a sense to numbers by which $1/2$ resembles 2, $1/3$ resembles 3, $1/4$ resembles 4, and so forth. To capture this sense, one can transform a function $f(z)$ into a function $f(w)$ by the change of variable³

$$w + \frac{1}{w} \leftarrow z,$$

or, more generally,

$$w + \frac{w_o^2}{w} \leftarrow z. \quad (10.1)$$

Equation (10.1) is *Vieta's transform*.⁴

For $|w| \gg |w_o|$, we have that $z \approx w$; but as $|w|$ approaches $|w_o|$ this ceases to be true. For $|w| \ll |w_o|$, $z \approx w_o^2/w$. The constant w_o is the *corner value*, in the neighborhood of which w transitions from the one domain to the other. Figure 10.1 plots Vieta's transform for real w in the case $w_o = 1$.

An interesting alternative to Vieta's transform is

$$w \parallel \frac{w_o^2}{w} \leftarrow z, \quad (10.2)$$

which in light of § 6.3 might be named *Vieta's parallel transform*.

Section 10.2 shows how Vieta's transform can be used.

10.2 Cubics

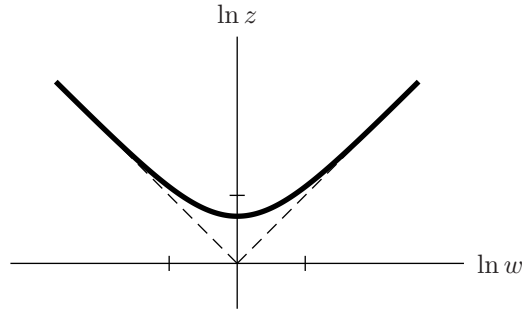
The general cubic polynomial is too hard to extract the roots of directly, so one begins by changing the variable

$$x + h \leftarrow z \quad (10.3)$$

²[45, "Cubic equation"] [45, "Quartic equation"] [46, "Quartic equation," 00:26, 9 Nov. 2006] [46, "François Viète," 05:17, 1 Nov. 2006] [46, "Gerolamo Cardano," 22:35, 31 Oct. 2006] [41, § 1.5]

³This change of variable broadly recalls the sum-of-exponentials form (5.19) of the $\cosh(\cdot)$ function, inasmuch as $\exp[-\phi] = 1/\exp \phi$.

⁴Also called "Vieta's substitution." [45, "Vieta's substitution"]

Figure 10.1: Vieta's transform (10.1) for $w_o = 1$, plotted logarithmically.

to obtain the polynomial

$$x^3 + (a_2 + 3h)x^2 + (a_1 + 2ha_2 + 3h^2)x + (a_0 + ha_1 + h^2a_2 + h^3).$$

The choice

$$h \equiv -\frac{a_2}{3} \quad (10.4)$$

casts the polynomial into the improved form

$$x^3 + \left[a_1 - \frac{a_2^2}{3} \right] x + \left[a_0 - \frac{a_1 a_2}{3} + 2 \left(\frac{a_2}{3} \right)^3 \right],$$

or better yet

$$x^3 - px - q,$$

where

$$\begin{aligned} p &\equiv -a_1 + \frac{a_2^2}{3}, \\ q &\equiv -a_0 + \frac{a_1 a_2}{3} - 2 \left(\frac{a_2}{3} \right)^3. \end{aligned} \quad (10.5)$$

The solutions to the equation

$$x^3 = px + q, \quad (10.6)$$

then, are the cubic polynomial's three roots.

So we have struck the $a_2 z^2$ term. That was the easy part; what to do next is not so obvious. If one could strike the px term as well, then the

roots would follow immediately, but no very simple substitution like (10.3) achieves this—or rather, such a substitution does achieve it, but at the price of reintroducing an unwanted x^2 or z^2 term. That way is no good. Lacking guidance, one might try many, various substitutions, none of which seems to help much; but after weeks or months of such frustration one might eventually discover Vieta's transform (10.1), with the idea of balancing the equation between offsetting w and $1/w$ terms. This works.

Vieta-transforming (10.6) by the change of variable

$$w + \frac{w_o^2}{w} \leftarrow x \quad (10.7)$$

we get the new equation

$$w^3 + (3w_o^2 - p)w + (3w_o^2 - p)\frac{w_o^2}{w} + \frac{w_o^6}{w^3} = q, \quad (10.8)$$

which invites the choice

$$w_o^2 \equiv \frac{p}{3}, \quad (10.9)$$

reducing (10.8) to read

$$w^3 + \frac{(p/3)^3}{w^3} = q.$$

Multiplying by w^3 and rearranging terms, we have the quadratic equation

$$(w^3)^2 = 2\left(\frac{q}{2}\right)w^3 - \left(\frac{p}{3}\right)^3, \quad (10.10)$$

which by (2.2) we know how to solve.

Vieta's transform has reduced the original cubic to a quadratic.

The careful reader will observe that (10.10) seems to imply six roots, double the three the fundamental theorem of algebra (§ 6.2.2) allows a cubic polynomial to have. We shall return to this point in § 10.3. For the moment, however, we should like to improve the notation by defining⁵

$$\begin{aligned} P &\leftarrow -\frac{p}{3}, \\ Q &\leftarrow +\frac{q}{2}, \end{aligned} \quad (10.11)$$

⁵Why did we not define P and Q so to begin with? Well, before unveiling (10.10), we lacked motivation to do so. To define inscrutable coefficients unnecessarily before the need for them is apparent seems poor applied mathematical style.

Table 10.1: The method to extract the three roots of the general cubic polynomial. (In the definition of w^3 , one can choose either sign.)

$$\begin{aligned}
0 &= z^3 + a_2 z^2 + a_1 z + a_0 \\
P &\equiv \frac{a_1}{3} - \left(\frac{a_2}{3}\right)^2 \\
Q &\equiv \frac{1}{2} \left[-a_0 + 3 \left(\frac{a_1}{3}\right) \left(\frac{a_2}{3}\right) - 2 \left(\frac{a_2}{3}\right)^3 \right] \\
w^3 &\equiv \begin{cases} 2Q & \text{if } P = 0, \\ Q \pm \sqrt{Q^2 + P^3} & \text{otherwise.} \end{cases} \\
x &\equiv \begin{cases} 0 & \text{if } P = 0 \text{ and } Q = 0, \\ w - P/w & \text{otherwise.} \end{cases} \\
z &= x - \frac{a_2}{3}
\end{aligned}$$

with which (10.6) and (10.10) are written

$$x^3 = 2Q - 3Px, \quad (10.12)$$

$$(w^3)^2 = 2Qw^3 + P^3. \quad (10.13)$$

Table 10.1 summarizes the complete cubic polynomial root extraction method in the revised notation—including a few fine points regarding superfluous roots and edge cases, treated in §§ 10.3 and 10.4 below.

10.3 Superfluous roots

As § 10.2 has observed, the equations of Table 10.1 seem to imply six roots, double the three the fundamental theorem of algebra (§ 6.2.2) allows a cubic polynomial to have. However, what the equations really imply is not six distinct roots but six distinct w . The definition $x \equiv w - P/w$ maps two w to any one x , so in fact the equations imply only three x and thus three roots z . The question then is: of the six w , which three do we really need and which three can we ignore as superfluous?

The six w naturally come in two groups of three: one group of three from the one w^3 and a second from the other. For this reason, we shall guess—and logically it is only a guess—that a single w^3 generates three distinct x and thus (because z differs from x only by a constant offset) all three roots z . If

the guess is right, then the second w^3 cannot but yield the same three roots, which means that the second w^3 is superfluous and can safely be overlooked. But is the guess right? Does a single w^3 in fact generate three distinct x ?

To prove that it does, let us suppose that it did not. Let us suppose that a single w^3 did generate two w which led to the same x . Letting the symbol w_1 represent the third w , then (since all three w come from the same w^3) the two w are $e^{+i2\pi/3}w_1$ and $e^{-i2\pi/3}w_1$. Because $x \equiv w - P/w$, by successive steps,

$$\begin{aligned} e^{+i2\pi/3}w_1 - \frac{P}{e^{+i2\pi/3}w_1} &= e^{-i2\pi/3}w_1 - \frac{P}{e^{-i2\pi/3}w_1}, \\ e^{+i2\pi/3}w_1 + \frac{P}{e^{-i2\pi/3}w_1} &= e^{-i2\pi/3}w_1 + \frac{P}{e^{+i2\pi/3}w_1}, \\ e^{+i2\pi/3} \left(w_1 + \frac{P}{w_1} \right) &= e^{-i2\pi/3} \left(w_1 + \frac{P}{w_1} \right), \end{aligned}$$

which can only be true if

$$w_1^2 = -P.$$

Cubing⁶ the last equation,

$$w_1^6 = -P^3;$$

but squaring the table's w^3 definition for $w = w_1$,

$$w_1^6 = 2Q^2 + P^3 \pm 2Q\sqrt{Q^2 + P^3}.$$

Combining the last two on w_1^6 ,

$$-P^3 = 2Q^2 + P^3 \pm 2Q\sqrt{Q^2 + P^3},$$

or, rearranging terms and halving,

$$Q^2 + P^3 = \mp Q\sqrt{Q^2 + P^3}.$$

Squaring,

$$Q^4 + 2Q^2P^3 + P^6 = Q^4 + Q^2P^3,$$

then canceling offsetting terms and factoring,

$$(P^3)(Q^2 + P^3) = 0.$$

⁶The verb *to cube* in this context means “to raise to the third power,” as to change y to y^3 , just as the verb *to square* means “to raise to the second power.”

The last equation demands rigidly that either $P = 0$ or $P^3 = -Q^2$. Some cubic polynomials do meet the demand—§ 10.4 will treat these and the reader is asked to set them aside for the moment—but most cubic polynomials do not meet it. For most cubic polynomials, then, the contradiction proves false the assumption which gave rise to it. The assumption: that the three x descending from a single w^3 were not distinct. Therefore, provided that $P \neq 0$ and $P^3 \neq -Q^2$, the three x descending from a single w^3 are indeed distinct, as was to be demonstrated.

The conclusion: *either, not both, of the two signs in the table's quadratic solution $w^3 \equiv Q \pm \sqrt{Q^2 + P^3}$ demands to be considered.* One can choose either sign; it matters not which.⁷ The one sign alone yields all three roots of the general cubic polynomial.

In calculating the three w from w^3 , one can apply the Newton-Raphson iteration (4.32), the Taylor series of Table 8.1, or any other convenient root-finding technique to find a single root w_1 such that $w_1^3 = w^3$. Then the other two roots come easier. They are $e^{\pm i2\pi/3}w_1$; but $e^{\pm i2\pi/3} = (-1 \pm i\sqrt{3})/2$, so

$$w = w_1, \frac{-1 \pm i\sqrt{3}}{2}w_1. \quad (10.14)$$

We should observe, incidentally, that nothing prevents two actual roots of a cubic polynomial from having the same value. This certainly is possible, and it does not mean that one of the two roots is superfluous or that the polynomial has fewer than three roots. For example, the cubic polynomial $(z - 1)(z - 1)(z - 2) = z^3 - 4z^2 + 5z - 2$ has roots at 1, 1 and 2, with a single root at $z = 2$ and a double root—that is, two roots—at $z = 1$. When this happens, the method of Table 10.1 properly yields the single root once and the double root twice, just as it ought to do.

10.4 Edge cases

Section 10.3 excepts the edge cases $P = 0$ and $P^3 = -Q^2$. Mostly the book does not worry much about edge cases, but the effects of these cubic edge cases seem sufficiently nonobvious that the book might include here a few words about them, if for no other reason than to offer the reader a model of how to think about edge cases on his own. Table 10.1 gives the quadratic solution

$$w^3 \equiv Q \pm \sqrt{Q^2 + P^3},$$

⁷Numerically, it can matter. As a simple rule, because w appears in the denominator of x 's definition, when the two w^3 differ in magnitude one might choose the larger.

in which § 10.3 generally finds it sufficient to consider either of the two signs. In the edge case $P = 0$,

$$w^3 = 2Q \text{ or } 0.$$

In the edge case $P^3 = -Q^2$,

$$w^3 = Q.$$

Both edge cases are interesting. In this section, we shall consider first the edge cases themselves, then their effect on the proof of § 10.3.

The edge case $P = 0$, like the general non-edge case, gives two distinct quadratic solutions w^3 . One of the two however is $w^3 = Q - Q = 0$, which is awkward in light of Table 10.1's definition that $x \equiv w - P/w$. For this reason, in applying the table's method when $P = 0$, one chooses the other quadratic solution, $w^3 = Q + Q = 2Q$.

The edge case $P^3 = -Q^2$ gives only the one quadratic solution $w^3 = Q$; or more precisely, it gives two quadratic solutions which happen to have the same value. This is fine. One merely accepts that $w^3 = Q$, and does not worry about choosing one w^3 over the other.

The double edge case, or *corner case*, arises where the two edges meet—where $P = 0$ and $P^3 = -Q^2$, or equivalently where $P = 0$ and $Q = 0$. At the corner, the trouble is that $w^3 = 0$ and that no alternate w^3 is available. However, according to (10.12), $x^3 = 2Q - 3Px$, which in this case means that $x^3 = 0$ and thus that $x = 0$ absolutely, no other x being possible. This implies the triple root $z = -a_2/3$.

Section 10.3 has excluded the edge cases from its proof of the sufficiency of a single w^3 . Let us now add the edge cases to the proof. In the edge case $P^3 = -Q^2$, both w^3 are the same, so the one w^3 suffices by default because the other w^3 brings nothing different. The edge case $P = 0$ however does give two distinct w^3 , one of which is $w^3 = 0$, which puts an awkward 0/0 in the table's definition of x . We address this edge in the spirit of l'Hôpital's rule, by sidestepping it, changing P infinitesimally from $P = 0$ to $P = \epsilon$. Then, choosing the $-$ sign in the definition of w^3 ,

$$\begin{aligned} w^3 &= Q - \sqrt{Q^2 + \epsilon^3} = Q - (Q) \left(1 + \frac{\epsilon^3}{2Q^2} \right) = -\frac{\epsilon^3}{2Q}, \\ w &= -\frac{\epsilon}{(2Q)^{1/3}}, \\ x &= w - \frac{\epsilon}{w} = -\frac{\epsilon}{(2Q)^{1/3}} + (2Q)^{1/3} = (2Q)^{1/3}. \end{aligned}$$

But choosing the + sign,

$$\begin{aligned} w^3 &= Q + \sqrt{Q^2 + \epsilon^3} = 2Q, \\ w &= (2Q)^{1/3}, \\ x &= w - \frac{\epsilon}{w} = (2Q)^{1/3} - \frac{\epsilon}{(2Q)^{1/3}} = (2Q)^{1/3}. \end{aligned}$$

Evidently the roots come out the same, either way. This completes the proof.

10.5 Quartics

Having successfully extracted the roots of the general cubic polynomial, we now turn our attention to the general quartic. The kernel of the cubic technique lay in reducing the cubic to a quadratic. The kernel of the quartic technique lies likewise in reducing the quartic to a cubic. The details differ, though; and, strangely enough, in some ways the quartic reduction is actually the simpler.⁸

As with the cubic, one begins solving the quartic by changing the variable

$$x + h \leftarrow z \tag{10.15}$$

to obtain the equation

$$x^4 = sx^2 + px + q, \tag{10.16}$$

where

$$\begin{aligned} h &\equiv -\frac{a_3}{4}, \\ s &\equiv -a_2 + 6\left(\frac{a_3}{4}\right)^2, \\ p &\equiv -a_1 + 2a_2\left(\frac{a_3}{4}\right) - 8\left(\frac{a_3}{4}\right)^3, \\ q &\equiv -a_0 + a_1\left(\frac{a_3}{4}\right) - a_2\left(\frac{a_3}{4}\right)^2 + 3\left(\frac{a_3}{4}\right)^4. \end{aligned} \tag{10.17}$$

⁸Even stranger, historically Ferrari discovered it earlier [45, “Quartic equation”]. Apparently Ferrari discovered the quartic’s resolvent cubic (10.22), which he could not solve until Tartaglia applied Vieta’s transform to it. What motivated Ferrari to chase the quartic solution while the cubic solution remained still unknown, this writer does not know, but one supposes that it might make an interesting story.

The reason the quartic is simpler to reduce is probably related to the fact that $(1)^{1/4} = \pm 1, \pm i$, whereas $(1)^{1/3} = 1, (-1 \pm i\sqrt{3})/2$. The $(1)^{1/4}$ brings a much neater result, the roots lying nicely along the Argand axes. This may also be why the quintic is intractable—but here we trespass the professional mathematician’s territory and stray from the scope of this book. See Ch. 6’s footnote 8.

To reduce (10.16) further, one must be cleverer. Ferrari⁹ supplies the cleverness. The clever idea is to transfer some but not all of the sx^2 term to the equation's left side by

$$x^4 + 2ux^2 = (2u + s)x^2 + px + q,$$

where u remains to be chosen; then to complete the square on the equation's left side as in § 2.2, but with respect to x^2 rather than x , as

$$(x^2 + u)^2 = k^2x^2 + px + j^2, \quad (10.18)$$

where

$$\begin{aligned} k^2 &\equiv 2u + s, \\ j^2 &\equiv u^2 + q. \end{aligned} \quad (10.19)$$

Now, one must regard (10.18) and (10.19) properly. In these equations, s , p and q have definite values fixed by (10.17), but not so u , j or k . The variable u is completely free; we have introduced it ourselves and can assign it any value we like. And though j^2 and k^2 depend on u , still, even after specifying u we remain free at least to choose signs for j and k . As for u , though no choice would truly be wrong, one supposes that a wise choice might at least render (10.18) easier to simplify.

So, what choice for u would be wise? Well, look at (10.18). The left side of that equation is a perfect square. The right side would be, too, if it were that $p = \pm 2jk$; so, arbitrarily choosing the $+$ sign, we propose the constraint that

$$p = 2jk, \quad (10.20)$$

or, better expressed,

$$j = \frac{p}{2k}. \quad (10.21)$$

Squaring (10.20) and substituting for j^2 and k^2 from (10.19), we have that

$$p^2 = 4(2u + s)(u^2 + q);$$

or, after distributing factors, rearranging terms and scaling, that

$$0 = u^3 + \frac{s}{2}u^2 + qu + \frac{4sq - p^2}{8}. \quad (10.22)$$

⁹[45, "Quartic equation"]

Equation (10.22) is the *resolvent cubic*, which we know by Table 10.1 how to solve for u , and which we now specify as a second constraint. If the constraints (10.21) and (10.22) are both honored, then we can safely substitute (10.20) into (10.18) to reach the form

$$(x^2 + u)^2 = k^2 x^2 + 2jkx + j^2,$$

which is

$$(x^2 + u)^2 = (kx + j)^2. \quad (10.23)$$

The resolvent cubic (10.22) of course yields three u not one, but the resolvent cubic is a voluntary constraint, so we can just pick one u and ignore the other two. Equation (10.19) then gives k (again, we can just pick one of the two signs), and (10.21) then gives j . With u , j and k established, (10.23) implies the quadratic

$$x^2 = \pm(kx + j) - u, \quad (10.24)$$

which (2.2) solves as

$$x = \pm \frac{k}{2} \pm_o \sqrt{\left(\frac{k}{2}\right)^2 \pm j - u}, \quad (10.25)$$

wherein the two \pm signs are tied together but the third, \pm_o sign is independent of the two. Equation (10.25), with the other equations and definitions of this section, reveals the four roots of the general quartic polynomial.

In view of (10.25), the change of variables

$$\begin{aligned} K &\leftarrow \frac{k}{2}, \\ J &\leftarrow j, \end{aligned} \quad (10.26)$$

improves the notation. Using the improved notation, Table 10.2 summarizes the complete quartic polynomial root extraction method.

10.6 Guessing the roots

It is entertaining to put pencil to paper and use Table 10.1's method to extract the roots of the cubic polynomial

$$0 = [z - 1][z - i][z + i] = z^3 - z^2 + z - 1.$$

Table 10.2: The method to extract the four roots of the general quartic polynomial. (In the table, the resolvent cubic is solved for u by the method of Table 10.1, where any one of the three resulting u serves. Either of the two K similarly serves. Of the three \pm signs in x 's definition, the \pm_o is independent but the other two are tied together, the four resulting combinations giving the four roots of the general quartic.)

$$\begin{aligned}
0 &= z^4 + a_3 z^3 + a_2 z^2 + a_1 z + a_0 \\
s &\equiv -a_2 + 6 \left(\frac{a_3}{4} \right)^2 \\
p &\equiv -a_1 + 2a_2 \left(\frac{a_3}{4} \right) - 8 \left(\frac{a_3}{4} \right)^3 \\
q &\equiv -a_0 + a_1 \left(\frac{a_3}{4} \right) - a_2 \left(\frac{a_3}{4} \right)^2 + 3 \left(\frac{a_3}{4} \right)^4 \\
0 &= u^3 + \frac{s}{2} u^2 + q u + \frac{4sq - p^2}{8} \\
K &\equiv \pm \frac{\sqrt{2u + s}}{2} \\
J &\equiv \begin{cases} \pm \sqrt{u^2 + q} & \text{if } K = 0, \\ p/4K & \text{otherwise.} \end{cases} \\
x &\equiv \pm K \pm_o \sqrt{K^2 \pm J - u} \\
z &= x - \frac{a_3}{4}
\end{aligned}$$

One finds that

$$\begin{aligned} z &= w + \frac{1}{3} - \frac{2}{3^2 w}, \\ w^3 &\equiv \frac{2(5 + \sqrt{3^3})}{3^3}, \end{aligned}$$

which says indeed that $z = 1, \pm i$, but just you try to simplify it! A more baroque, more impenetrable way to write the number 1 is not easy to conceive. One has found the number 1 but cannot recognize it. Figuring the square and cube roots in the expression numerically, the root of the polynomial comes mysteriously to 1.0000, but why? The root's symbolic form gives little clue.

In general no better way is known;¹⁰ we are stuck with the cubic baroquity. However, to the extent that a cubic, a quartic, a quintic or any other polynomial has real, rational roots, a trick is known to sidestep Tables 10.1 and 10.2 and guess the roots directly. Consider for example the quintic polynomial

$$z^5 - \frac{7}{2}z^4 + 4z^3 + \frac{1}{2}z^2 - 5z + 3.$$

Doubling to make the coefficients all integers produces the polynomial

$$2z^5 - 7z^4 + 8z^3 + 1z^2 - 10Az + 6,$$

which naturally has the same roots. If the roots are complex or irrational, they are hard to guess; but if any of the roots happens to be real and rational, it must belong to the set

$$\left\{ \pm 1, \pm 2, \pm 3, \pm 6, \pm \frac{1}{2}, \pm \frac{2}{2}, \pm \frac{3}{2}, \pm \frac{6}{2} \right\}.$$

No other real, rational root is possible. Trying the several candidates on the polynomial, one finds that 1, -1 and 3/2 are indeed roots. Dividing these out leaves a quadratic which is easy to solve for the remaining roots.

The real, rational candidates are the factors of the polynomial's trailing coefficient (in the example, 6, whose factors are $\pm 1, \pm 2, \pm 3$ and ± 6) divided by the factors of the polynomial's leading coefficient (in the example, 2, whose factors are ± 1 and ± 2). The reason no other real, rational root is

¹⁰At least, no better way is known to this author. If any reader can straightforwardly simplify the expression without solving a cubic polynomial of some kind, the author would like to hear of it.

possible is seen¹¹ by writing $z = p/q$ —where p and q are integers and the fraction p/q is fully reduced—then multiplying the n th-order polynomial by q^n to reach the form

$$a_n p^n + a_{n-1} p^{n-1} q + \cdots + a_1 p q^{n-1} + a_0 q^n = 0,$$

where all the coefficients a_k are integers. Moving the q^n term to the equation's right side, we have that

$$(a_n p^{n-1} + a_{n-1} p^{n-2} q + \cdots + a_1 q^{n-1}) p = -a_0 q^n,$$

which implies that $a_0 q^n$ is a multiple of p . But by demanding that the fraction p/q be fully reduced, we have defined p and q to be *relatively prime* to one another—that is, we have defined them to have no factors but ± 1 in common—so, not only $a_0 q^n$ but a_0 itself is a multiple of p . By similar reasoning, a_n is a multiple of q . But if a_0 is a multiple of p , and a_n , a multiple of q , then p and q are factors of a_0 and a_n respectively. We conclude for this reason, as was to be demonstrated, that no real, rational root is possible except a factor of a_0 divided by a factor of a_n .¹²

Such root-guessing is little more than an algebraic trick, of course, but it can be a pretty useful trick if it saves us the embarrassment of inadvertently expressing simple rational numbers in ridiculous ways.

One could write much more about higher-order algebra, but now that the reader has tasted the topic he may feel inclined to agree that, though the general methods this chapter has presented to solve cubics and quartics are interesting, further effort were nevertheless probably better spent elsewhere. The next several chapters turn to the topic of the matrix, harder but much more profitable, toward which we mean to put substantial effort.

¹¹The presentation here is quite informal. We do not want to spend many pages on this.

¹²[41, § 3.2]

Chapter 11

The matrix

Chapters 2 through 9 have laid solidly the basic foundations of applied mathematics. This chapter begins to build on those foundations, demanding some heavier mathematical lifting.

Taken by themselves, most of the foundational methods of the earlier chapters have handled only one or at most a few numbers (or functions) at a time. However, in practical applications the need to handle large arrays of numbers at once arises often. Some nonobvious effects emerge then, as, for example, the eigenvalue of Ch. 14.

Regarding the eigenvalue: the eigenvalue was always there, but prior to this point in the book it was usually trivial—the eigenvalue of 5 is just 5, for instance—so we didn’t bother much to talk about it. It is when numbers are laid out in orderly grids like

$$C = \begin{bmatrix} 6 & 4 & 0 \\ 3 & 0 & 1 \\ 3 & 1 & 0 \end{bmatrix}$$

that nontrivial eigenvalues arise (though you cannot tell just by looking, the eigenvalues of C happen to be -1 and $[7 \pm \sqrt{0x49}]/2$). But, just what is an *eigenvalue*? Answer: an eigenvalue is the value by which an object like C scales an eigenvector without altering the eigenvector’s direction. Of course, we have not yet said what an *eigenvector* is, either, or how C might scale something, but it is to answer precisely such questions that this chapter and the three which follow it are written.

So, we are getting ahead of ourselves. Let’s back up.

An object like C is called a *matrix*. It serves as a generalized coefficient or multiplier. Where we have used single numbers as coefficients or multipliers

heretofore, one can with sufficient care often use matrices instead. The matrix interests us for this reason among others.

The technical name for the “single number” is the *scalar*. Such a number, as for instance 5 or $-4 + i3$, is called a scalar because its action alone during multiplication is simply to scale the thing it multiplies. Besides acting alone, however, scalars can also act in concert—in orderly formations—thus constituting any of three basic kinds of arithmetical object:

- the *scalar* itself, a single number like $\alpha = 5$ or $\beta = -4 + i3$;
- the *vector*, a column of m scalars like

$$\mathbf{u} = \begin{bmatrix} 5 \\ -4 + i3 \end{bmatrix},$$

which can be written in-line with the notation $\mathbf{u} = [5 \ -4 + i3]^T$ (here there are two scalar elements, 5 and $-4 + i3$, so in this example $m = 2$);

- the *matrix*, an $m \times n$ grid of scalars, or equivalently a row of n vectors, like

$$A = \begin{bmatrix} 0 & 6 & 2 \\ 1 & 1 & -1 \end{bmatrix},$$

which can be written in-line with the notation $A = [0 \ 6 \ 2; 1 \ 1 \ -1]$ or the notation $A = [0 \ 1; 6 \ 1; 2 \ -1]^T$ (here there are two rows and three columns of scalar elements, so in this example $m = 2$ and $n = 3$).

Several general points are immediately to be observed about these various objects. First, despite the geometrical Argand interpretation of the complex number, a complex number is not a two-element vector but a scalar; therefore any or all of a vector’s or matrix’s scalar elements can be complex. Second, an m -element vector does not differ for most purposes from an $m \times 1$ matrix; generally the two can be regarded as the same thing. Third, the three-element (that is, three-dimensional) geometrical vector of § 3.3 is just an m -element vector with $m = 3$. Fourth, m and n can be any nonnegative integers, even one, even zero, even infinity.¹

Where one needs visually to distinguish a symbol like A representing a matrix, one can write it $[A]$, in square brackets.² Normally however a simple A suffices.

¹Fifth, though the progression *scalar*, *vector*, *matrix* suggests next a “matrix stack” or stack of p matrices, such objects in fact are seldom used. As we shall see in § 11.1, the chief advantage of the standard matrix is that it neatly represents the linear transformation of one vector into another. “Matrix stacks” bring no such advantage. This book does not treat them.

²Alternate notations seen in print include $\overline{\overline{A}}$ and \mathbf{A} .

The matrix is a notoriously hard topic to motivate. The idea of the matrix is deceptively simple. The mechanics of matrix arithmetic are deceptively intricate. The most basic body of matrix theory, without which little or no useful matrix work can be done, is deceptively extensive. The matrix neatly encapsulates a substantial knot of arithmetical tedium and clutter, but to understand the matrix one must first understand the tedium and clutter the matrix encapsulates. As far as the author is aware, no one has ever devised a way to introduce the matrix which does not seem shallow, tiresome, irksome, even interminable at first encounter; yet the matrix is too important to ignore. Applied mathematics brings nothing else quite like it.³

Chapters 11 through 14 treat the matrix and its algebra. This chapter,

³In most of its chapters, the book seeks a balance between terseness the determined beginner cannot penetrate and prolixity the seasoned veteran will not abide. The matrix upsets this balance.

Part of the trouble with the matrix is that its arithmetic is just that, an arithmetic, no more likely to be mastered by mere theoretical study than was the classical arithmetic of childhood. To master matrix arithmetic, one must drill it; yet the book you hold is fundamentally one of theory not drill.

The reader who has previously drilled matrix arithmetic will meet here the essential applied theory of the matrix. That reader will find this chapter and the next three tedious enough. The reader who has not previously drilled matrix arithmetic, however, is likely to find these chapters positively hostile. Only the doggedly determined beginner will learn the matrix here alone; others will find it more amenable to drill matrix arithmetic first in the early chapters of an introductory linear algebra textbook, dull though such chapters be (see [27] or better yet the fine, surprisingly less dull [20] for instance, though the early chapters of almost any such book give the needed arithmetical drill.) Returning here thereafter, the beginner can expect to find *these* chapters still tedious but no longer impenetrable. The reward is worth the effort. That is the approach the author recommends.

To the mathematical rebel, the young warrior with face painted and sword agleam, still determined to learn the matrix here alone, the author salutes his honorable defiance. Would the rebel consider alternate counsel? If so, then the rebel might compose a dozen matrices of various sizes and shapes, broad, square and tall, decomposing each carefully by pencil per the Gauss-Jordan method of § 12.3, checking results (again by pencil; using a machine defeats the point of the exercise, and using a sword, well, it won't work) by multiplying factors to restore the original matrices. Several hours of such drill should build the young warrior the practical arithmetical foundation to master—with commensurate effort—the theory these chapters bring. The way of the warrior is hard, but conquest is not impossible.

To the matrix veteran, the author presents these four chapters with grim enthusiasm. Substantial, logical, necessary the chapters may be, but exciting they are not. At least, the earlier parts are not very exciting (later parts are better). As a reasonable compromise, the veteran seeking more interesting reading might skip directly to Chs. 13 and 14, referring back to Chs. 11 and 12 as need arises.

Ch. 11, introduces the rudiments of the matrix itself.⁴

11.1 Provenance and basic use

It is in the study of linear transformations that the concept of the matrix first arises. We begin there.

11.1.1 The linear transformation

Section 7.3.3 has introduced the idea of linearity. The *linear transformation*⁵ is the operation of an $m \times n$ matrix A , as in

$$A\mathbf{x} = \mathbf{b}, \quad (11.1)$$

to transform an n -element vector \mathbf{x} into an m -element vector \mathbf{b} , while respecting the rules of linearity

$$\begin{aligned} A(\mathbf{x}_1 + \mathbf{x}_2) &= A\mathbf{x}_1 + A\mathbf{x}_2 = \mathbf{b}_1 + \mathbf{b}_2, \\ A(\alpha\mathbf{x}) &= \alpha A\mathbf{x} = \alpha\mathbf{b}, \\ A(0) &= 0. \end{aligned} \quad (11.2)$$

For example,

$$A = \begin{bmatrix} 0 & 6 & 2 \\ 1 & 1 & -1 \end{bmatrix}$$

is the 2×3 matrix which transforms a three-element vector \mathbf{x} into a two-element vector \mathbf{b} such that

$$A\mathbf{x} = \begin{bmatrix} 0x_1 + 6x_2 + 2x_3 \\ 1x_1 + 1x_2 - 1x_3 \end{bmatrix} = \mathbf{b},$$

where

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}.$$

⁴[2][13][20][27]

⁵Professional mathematicians conventionally are careful to begin by drawing a clear distinction between the ideas of the linear transformation, the basis set and the simultaneous system of linear equations—proving from suitable axioms that the three amount more or less to the same thing, rather than implicitly assuming the fact. The professional approach [2, Chs. 1 and 2][27, Chs. 1, 2 and 5] has much to recommend it, but it is not the approach we will follow here.

In general, the operation of a matrix A is that^{6,7}

$$b_i = \sum_{j=1}^n a_{ij}x_j, \quad (11.3)$$

where x_j is the j th element of \mathbf{x} , b_i is the i th element of \mathbf{b} , and

$$a_{ij} \equiv [A]_{ij}$$

is the element at the i th row and j th column of A , counting from top left (in the example for instance, $a_{12} = 6$).

Besides representing linear transformations as such, matrices can also represent simultaneous systems of linear equations. For example, the system

$$\begin{aligned} 0x_1 + 6x_2 + 2x_3 &= 2, \\ 1x_1 + 1x_2 - 1x_3 &= 4, \end{aligned}$$

is compactly represented as

$$A\mathbf{x} = \mathbf{b},$$

with A as given above and $\mathbf{b} = [2 \ 4]^T$. Seen from this point of view, a simultaneous system of linear equations is itself neither more nor less than a linear transformation.

11.1.2 Matrix multiplication (and addition)

Nothing prevents one from lining several vectors \mathbf{x}_k up in a row, industrial mass production-style, transforming them at once into the corresponding

⁶As observed in Appendix B, there are unfortunately not enough distinct Roman and Greek letters available to serve the needs of higher mathematics. In matrix work, the Roman letters ijk conventionally serve as indices, but the same letter i also serves as the imaginary unit, which is not an index and has nothing to do with indices. Fortunately, the meaning is usually clear from the context: i in \sum_i or a_{ij} is an index; i in $-4 + i3$ or $e^{i\phi}$ is the imaginary unit. Should a case arise in which the meaning is not clear, one can use ℓjk or some other convenient letters for the indices.

⁷Whether to let the index j run from 0 to $n-1$ or from 1 to n is an awkward question of applied mathematical style. In computers, the index normally runs from 0 to $n-1$, and in many ways this really is the more sensible way to do it. In mathematical theory, however, a 0 index normally implies something special or basic about the object it identifies. The book you are reading tends to let the index run from 1 to n , following mathematical convention in the matter for this reason.

Conceived more generally, an $m \times n$ matrix can be considered an $\infty \times \infty$ matrix with zeros in the unused cells. Here, both indices i and j run from $-\infty$ to $+\infty$ anyway, so the computer's indexing convention poses no dilemma in this case. See § 11.3.

vectors \mathbf{b}_k by the same matrix A . In this case,

$$\begin{aligned} X &\equiv [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_p], \\ B &\equiv [\mathbf{b}_1 \ \mathbf{b}_2 \ \cdots \ \mathbf{b}_p], \\ AX &= B, \\ b_{ik} &= \sum_{j=1}^n a_{ij}x_{jk}. \end{aligned} \tag{11.4}$$

Equation (11.4) implies a definition for matrix multiplication. Such matrix multiplication is associative since

$$\begin{aligned} [(A)(XY)]_{ik} &= \sum_{j=1}^n a_{ij}[XY]_{jk} \\ &= \sum_{j=1}^n a_{ij} \left[\sum_{\ell=1}^p x_{j\ell}y_{\ell k} \right] \\ &= \sum_{\ell=1}^p \sum_{j=1}^n a_{ij}x_{j\ell}y_{\ell k} \\ &= [(AX)(Y)]_{ik}. \end{aligned} \tag{11.5}$$

Matrix multiplication is not generally commutative, however;

$$AX \neq XA, \tag{11.6}$$

as one can show by a suitable counterexample like $A = [0 \ 1; 0 \ 0]$, $X = [1 \ 0; 0 \ 0]$. To multiply a matrix by a scalar, one multiplies each of the matrix's elements individually by the scalar:

$$[\alpha A]_{ij} = \alpha a_{ij}. \tag{11.7}$$

Evidently multiplication by a scalar is commutative: $\alpha A\mathbf{x} = A\alpha\mathbf{x}$.

Matrix addition works in the way one would expect, element by element; and as one can see from (11.4), under multiplication, matrix addition is indeed distributive:

$$\begin{aligned} [X + Y]_{ij} &= x_{ij} + y_{ij}; \\ (A)(X + Y) &= AX + AY; \\ (A + C)(X) &= AX + CX. \end{aligned} \tag{11.8}$$

11.1.3 Row and column operators

The matrix equation $A\mathbf{x} = \mathbf{b}$ represents the linear transformation of \mathbf{x} into \mathbf{b} , as we have seen. Viewed from another perspective, however, the same matrix equation represents something else; it represents a weighted sum of the columns of A , with the elements of \mathbf{x} as the weights. In this view, one writes (11.3) as

$$\mathbf{b} = \sum_{j=1}^n [A]_{*j} x_j, \quad (11.9)$$

where $[A]_{*j}$ is the j th column of A . Here \mathbf{x} is not only a vector; it is also an operator. It operates on A 's columns. By virtue of multiplying A from the right, the vector \mathbf{x} is a *column operator* acting on A .

If several vectors \mathbf{x}_k line up in a row to form a matrix X , such that $AX = B$, then the matrix X is likewise a column operator:

$$[B]_{*k} = \sum_{j=1}^n [A]_{*j} x_{jk}. \quad (11.10)$$

The k th column of X weights the several columns of A to yield the k th column of B .

If a matrix multiplying from the right is a column operator, is a matrix multiplying from the left a *row operator*? Indeed it is. Another way to write $AX = B$, besides (11.10), is

$$[B]_{i*} = \sum_{j=1}^n a_{ij} [X]_{j*}. \quad (11.11)$$

The i th row of A weights the several rows of X to yield the i th row of B . The matrix A is a row operator. (Observe the notation. The $*$ here means “any” or “all.” Hence $[X]_{j*}$ means “ j th row, all columns of X ”—that is, the j th row of X . Similarly, $[A]_{*j}$ means “all rows, j th column of A ”—that is, the j th column of A .)

Column operators attack from the right; row operators, from the left. This rule is worth memorizing; the concept is important. In $AX = B$, the matrix X operates on A 's columns; the matrix A operates on X 's rows.

Since matrix multiplication produces the same result whether one views it as a linear transformation (11.4), a column operation (11.10) or a row operation (11.11), one might wonder what purpose lies in defining matrix multiplication three separate ways. However, it is not so much for the sake

of the mathematics that we define it three ways as it is for the sake of the mathematician. We do it for ourselves. Mathematically, the latter two do indeed expand to yield (11.4), but as written the three represent three different perspectives on the matrix. A tedious, nonintuitive matrix theorem from one perspective can appear suddenly obvious from another (see for example eqn. 11.63). Results hard to visualize one way are easy to visualize another. It is worth developing the mental agility to view and handle matrices all three ways for this reason.

11.1.4 The transpose and the adjoint

One function peculiar to matrix algebra is the *transpose*

$$\begin{aligned} C &= A^T, \\ c_{ij} &= a_{ji}, \end{aligned} \tag{11.12}$$

which mirrors an $m \times n$ matrix into an $n \times m$ matrix. For example,

$$A^T = \begin{bmatrix} 0 & 1 \\ 6 & 1 \\ 2 & -1 \end{bmatrix}.$$

Similar and even more useful is the *conjugate transpose* or *adjoint*⁸

$$\begin{aligned} C &= A^*, \\ c_{ij} &= a_{ji}^*, \end{aligned} \tag{11.13}$$

which again mirrors an $m \times n$ matrix into an $n \times m$ matrix, but conjugates each element as it goes.

The transpose is convenient notationally to write vectors and matrices in-line and to express certain matrix-arithmetical mechanics; but algebraically the transpose is artificial. It is the adjoint rather which mirrors a matrix properly. (If the transpose and adjoint functions applied to words as to matrices, then the transpose of “derivations” would be “snoitavired,” whereas the adjoint would be “snoitavirəb.” See the difference?) On real-valued matrices like the A in the example, of course, the transpose and the adjoint amount to the same thing.

⁸Alternate notations sometimes seen in print for the adjoint include A^\dagger (a notation which in this book means something unrelated) and A^H (a notation which recalls the name of the mathematician Charles Hermite). However, the book you are reading writes the adjoint only as A^* , a notation which better captures the sense of the thing in the author’s view.

If one needed to conjugate the elements of a matrix without transposing the matrix itself, one could contrive notation like A^{*T} . Such a need seldom arises, however.

Observe that

$$\begin{aligned}(A_2 A_1)^T &= A_1^T A_2^T, \\ (A_2 A_1)^* &= A_1^* A_2^*,\end{aligned}\tag{11.14}$$

and more generally that⁹

$$\begin{aligned}\left(\prod_k A_k\right)^T &= \prod_k A_k^T, \\ \left(\prod_k A_k\right)^* &= \prod_k A_k^*.\end{aligned}\tag{11.15}$$

11.2 The Kronecker delta

Section 7.7 has introduced the Dirac delta. The discrete analog of the Dirac delta is the *Kronecker delta*¹⁰

$$\delta_i \equiv \begin{cases} 1 & \text{if } i = 0, \\ 0 & \text{otherwise;} \end{cases}\tag{11.16}$$

or

$$\delta_{ij} \equiv \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{otherwise.} \end{cases}\tag{11.17}$$

The Kronecker delta enjoys the Dirac-like properties that

$$\sum_{i=-\infty}^{\infty} \delta_i = \sum_{i=-\infty}^{\infty} \delta_{ij} = \sum_{j=-\infty}^{\infty} \delta_{ij} = 1\tag{11.18}$$

and that

$$\sum_{j=-\infty}^{\infty} \delta_{ij} a_{jk} = a_{ik},\tag{11.19}$$

the latter of which is the Kronecker sifting property. The Kronecker equations (11.18) and (11.19) parallel the Dirac equations (7.13) and (7.14).

⁹Recall from § 2.3 that $\prod_k A_k = \cdots A_3 A_2 A_1$, whereas $\coprod_k A_k = A_1 A_2 A_3 \cdots$.

¹⁰[46, “Kronecker delta,” 15:59, 31 May 2006]

11.3 Dimensionality and matrix forms

An $m \times n$ matrix like

$$X = \begin{bmatrix} -4 & 0 \\ 1 & 2 \\ 2 & -1 \end{bmatrix}$$

can be viewed as the $\infty \times \infty$ matrix

$$X = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & -4 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 2 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 2 & -1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

with zeros in the unused cells. As before, $x_{11} = -4$ and $x_{32} = -1$, but now x_{ij} exists for all integral i and j ; for instance, $x_{(-1)(-1)} = 0$. For such a matrix, indeed for all matrices, the matrix multiplication rule (11.4) generalizes to

$$\begin{aligned} B &= AX, \\ b_{ik} &= \sum_{j=-\infty}^{\infty} a_{ij}x_{jk}. \end{aligned} \tag{11.20}$$

For square matrices whose purpose is to manipulate other matrices or vectors in place, merely padding with zeros often does not suit. Consider for example the square matrix

$$A_3 = \begin{bmatrix} 1 & 0 & 0 \\ 5 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

This A_3 is indeed a matrix, but when it acts A_3X as a row operator on some $3 \times p$ matrix X , its effect is to add to X 's second row, 5 times the first. Further consider

$$A_4 = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 5 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

which does the same to a $4 \times p$ matrix X . We can also define A_5, A_6, A_7, \dots , if we want; but, really, all these express the same operation: “to add to the second row, 5 times the first.”

The $\infty \times \infty$ matrix

$$A = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 5 & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \end{bmatrix}$$

expresses the operation generally. As before, $a_{11} = 1$ and $a_{21} = 5$, but now also $a_{(-1)(-1)} = 1$ and $a_{09} = 0$, among others. By running ones infinitely both ways out the *main diagonal*, we guarantee by (11.20) that when A acts AX on a matrix X of any dimensionality whatsoever, A adds to the second row of X , 5 times the first—and affects no other row. (But what if X is a $1 \times p$ matrix, and *has* no second row? Then the operation AX creates a new second row, 5 times the first—or rather so fills in X 's previously null second row.)

In the infinite-dimensional view, the matrices A and X differ essentially.¹¹ This section explains, developing some nonstandard formalisms the derivations of later sections and chapters can use.¹²

¹¹This particular section happens to use the symbols A and X to represent certain specific matrix forms because such usage flows naturally from the usage $A\mathbf{x} = \mathbf{b}$ of § 11.1. Such usage admittedly proves awkward in other contexts. Traditionally in matrix work and elsewhere in the book, the letter A does not necessarily represent an extended operator as it does here, but rather an arbitrary matrix of no particular form.

¹²The idea of infinite dimensionality is sure to discomfit some readers, who have studied matrices before and are used to thinking of a matrix as having some definite size. There is nothing wrong with thinking of a matrix as having some definite size, only that that view does not suit the present book's development. And really, the idea of an $\infty \times 1$ vector or an $\infty \times \infty$ matrix should not seem so strange. After all, consider the vector \mathbf{u} such that

$$u_\ell = \sin \ell \epsilon,$$

where $0 < \epsilon \ll 1$ and ℓ is an integer, which holds all values of the function $\sin \theta$ of a real argument θ . Of course one does not actually write down or store all the elements of an infinite-dimensional vector or matrix, any more than one actually writes down or stores all the bits (or digits) of 2π . Writing them down or storing them is not the point. The point is that infinite dimensionality is all right; that the idea thereof does not threaten to overturn the reader's preëxisting matrix knowledge; that, though the construct seem unfamiliar, no fundamental conceptual barrier rises against it.

Different ways of looking at the same mathematics can be extremely useful to the applied mathematician. The applied mathematical reader who has never heretofore considered

11.3.1 The null and dimension-limited matrices

The *null matrix* is just what its name implies:

$$0 = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix};$$

or more compactly,

$$[0]_{ij} = 0.$$

Special symbols like $\bar{0}$, $\mathbf{0}$ or O are possible for the null matrix, but usually a simple 0 suffices. There are no surprises here; the null matrix brings all the expected properties of a zero, like

$$\begin{aligned} 0 + A &= A, \\ [0][X] &= 0. \end{aligned}$$

The same symbol 0 used for the null scalar (zero) and the null matrix is used for the null vector, too. Whether the scalar 0, the vector 0 and the matrix 0 actually represent different things is a matter of semantics, but the three are interchangeable for most practical purposes in any case. Basically, a zero is a zero; there's not much else to it.¹³

Now a formality: the ordinary $m \times n$ matrix X can be viewed, infinite-dimensionally, as a variation on the null matrix, inasmuch as X differs from the null matrix only in the mn elements x_{ij} , $1 \leq i \leq m$, $1 \leq j \leq n$. Though the theoretical dimensionality of X be $\infty \times \infty$, one need record only the mn elements, plus the values of m and n , to retain complete information about such a matrix. So the semantics are these: when we call a matrix X an $m \times n$ matrix, or more precisely a *dimension-limited matrix* with an $m \times n$

infinite dimensionality in vectors and matrices would be well served to take the opportunity to do so here. As we shall discover in Ch. 12, dimensionality is a poor measure of a matrix's size in any case. What really counts is not a matrix's $m \times n$ dimensionality but rather its *rank*.

¹³Well, of course, there's a lot else to it, when it comes to dividing by zero as in Ch. 4, or to summing an infinity of zeros as in Ch. 7, but those aren't what we were speaking of here.

active region, we shall mean formally that X is an $\infty \times \infty$ matrix whose elements are all zero outside the $m \times n$ rectangle:

$$x_{ij} = 0 \text{ except where } 1 \leq i \leq m \text{ and } 1 \leq j \leq n. \quad (11.21)$$

By these semantics, every 3×2 matrix (for example) is also a formally a 4×4 matrix; but a 4×4 matrix is not in general a 3×2 matrix.

11.3.2 The identity and scalar matrices and the extended operator

The *general identity matrix*—or simply, the *identity matrix*—is

$$I = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & 0 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \end{bmatrix},$$

or more compactly,

$$[I]_{ij} = \delta_{ij}, \quad (11.22)$$

where δ_{ij} is the Kronecker delta of § 11.2. The identity matrix I is a matrix 1, as it were,¹⁴ bringing the essential property one expects of a 1:

$$IX = X = XI. \quad (11.23)$$

The *scalar matrix* is

$$\lambda I = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & \lambda & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & \lambda & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & \lambda & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & \lambda & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & \lambda & 0 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \end{bmatrix},$$

or more compactly,

$$[\lambda I]_{ij} = \lambda \delta_{ij}, \quad (11.24)$$

¹⁴In fact you can write it as 1 if you like. That is essentially what it is. The I can be regarded as standing for “identity” or as the Roman numeral I.

If the identity matrix I is a matrix 1, then the scalar matrix λI is a matrix λ , such that

$$[\lambda I]X = \lambda X = X[\lambda I]. \quad (11.25)$$

The identity matrix is (to state the obvious) just the scalar matrix with $\lambda = 1$.

The *extended operator* A is a variation on the scalar matrix λI , $\lambda \neq 0$, inasmuch as A differs from λI only in p specific elements, with p a finite number. Symbolically,

$$a_{ij} = \begin{cases} (\lambda)(\delta_{ij} + \alpha_k) & \text{if } (i, j) = (i_k, j_k), 1 \leq k \leq p, \\ \lambda \delta_{ij} & \text{otherwise;} \end{cases} \quad (11.26)$$

$$\lambda \neq 0.$$

The several α_k control how the extended operator A differs from λI . One need record only the several α_k along with their respective addresses (i_k, j_k) , plus the scale λ , to retain complete information about such a matrix. For example, for an extended operator fitting the pattern

$$A = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & \lambda & 0 & 0 & 0 & 0 & \cdots \\ \cdots & \lambda \alpha_1 & \lambda & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & \lambda & \lambda \alpha_2 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & \lambda(1 + \alpha_3) & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & \lambda & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

one need record only the values of α_1 , α_2 and α_3 , the respective addresses $(2, 1)$, $(3, 4)$ and $(4, 4)$, and the value of the scale λ ; this information alone implies the entire $\infty \times \infty$ matrix A .

When we call a matrix A an *extended $n \times n$ operator*, or an extended operator with an $n \times n$ *active region*, we shall mean formally that A is an $\infty \times \infty$ matrix and is further an extended operator for which

$$1 \leq i_k \leq n \text{ and } 1 \leq j_k \leq n \text{ for all } 1 \leq k \leq p. \quad (11.27)$$

That is, an extended $n \times n$ operator is one whose several α_k all lie within the $n \times n$ square. The A in the example is an extended 4×4 operator (and also a 5×5 , a 6×6 , etc., but not a 3×3).

(Often in practice for smaller operators—especially in the typical case that $\lambda = 1$ —one finds it easier just to record all the $n \times n$ elements of the active region. This is fine. Large matrix operators however tend to be

sparse, meaning that they depart from λI in only a very few of their many elements. It would waste a lot of computer memory explicitly to store all those zeros, so one normally stores just the few elements, instead.)

Implicit in the definition of the extended operator is that the identity matrix I and the scalar matrix λI , $\lambda \neq 0$, are extended operators with 0×0 active regions (and also 1×1 , 2×2 , etc.). If $\lambda = 0$, however, the scalar matrix λI is just the null matrix, which is no extended operator but rather by definition a 0×0 dimension-limited matrix.

11.3.3 The active region

Though maybe obvious, it bears stating explicitly that a product of dimension-limited and/or extended-operational matrices with $n \times n$ active regions itself has an $n \times n$ active region.¹⁵ (Remember that a matrix with an $m' \times n'$ active region also by definition has an $n \times n$ active region if $m' \leq n$ and $n' \leq n$.) If any of the factors has dimension-limited form then so does the product; otherwise the product is an extended operator.¹⁶

11.3.4 Other matrix forms

Besides the dimension-limited form of § 11.3.1 and the extended-operational form of § 11.3.2, other infinite-dimensional matrix forms are certainly possible. One could for example advantageously define a “null sparse” form, recording only nonzero elements and their addresses in an otherwise null matrix; or a “tridiagonal extended” form, bearing repeated entries not only along the main diagonal but also along the diagonals just above and just below. Section 11.9 introduces one worthwhile matrix which fits neither the dimension-limited nor the extended-operational form. Still, the dimension-

¹⁵The section’s earlier subsections formally define the term *active region* with respect to each of the two matrix forms.

¹⁶If symbolic proof of the subsection’s claims is wanted, here it is in outline:

$$\begin{aligned}
 a_{ij} &= \lambda_a \delta_{ij} && \text{unless } 1 \leq (i, j) \leq n, \\
 b_{ij} &= \lambda_b \delta_{ij} && \text{unless } 1 \leq (i, j) \leq n; \\
 [AB]_{ij} &= \sum_k a_{ik} b_{kj} \\
 &= \begin{cases} \sum_k (\lambda_a \delta_{ik}) b_{kj} = \lambda_a b_{ij} & \text{unless } 1 \leq i \leq n \\ \sum_k a_{ik} (\lambda_b \delta_{kj}) = \lambda_b a_{ij} & \text{unless } 1 \leq j \leq n \end{cases} \\
 &= \lambda_a \lambda_b \delta_{ij} && \text{unless } 1 \leq (i, j) \leq n.
 \end{aligned}$$

It’s probably easier just to sketch the matrices and look at them, though.

limited and extended-operational forms are normally the most useful, and they are the ones we shall principally be handling in this book.

One reason to have defined specific infinite-dimensional matrix forms is to show how straightforwardly one can fully represent a practical matrix of an infinity of elements by a modest, finite quantity of information. Further reasons to have defined such forms will soon occur.

11.3.5 The rank- r identity matrix

The *rank- r identity matrix* I_r is the dimension-limited matrix for which

$$[I_r]_{ij} = \begin{cases} \delta_{ij} & \text{if } 1 \leq i \leq r \text{ and/or } 1 \leq j \leq r, \\ 0 & \text{otherwise,} \end{cases} \quad (11.28)$$

where either the “and” or the “or” can be regarded (it makes no difference). The effect of I_r is that

$$\begin{aligned} I_m X &= X = X I_n, \\ I_m \mathbf{x} &= \mathbf{x}, \end{aligned} \quad (11.29)$$

where X is an $m \times n$ matrix and \mathbf{x} , an $m \times 1$ vector. Examples of I_r include

$$I_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

(Remember that in the infinite-dimensional view, I_3 , though a 3×3 matrix, is formally an $\infty \times \infty$ matrix with zeros in the unused cells. It has only the three ones and fits the 3×3 dimension-limited form of § 11.3.1. The areas of I_3 not shown are all zero, even along the main diagonal.)

The rank r can be any nonnegative integer, even zero (though the rank-zero identity matrix I_0 is in fact the null matrix, normally just written 0). If alternate indexing limits are needed (for instance for a computer-indexed identity matrix whose indices run from 0 to $r - 1$), the notation I_a^b , where

$$[I_a^b]_{ij} \equiv \begin{cases} \delta_{ij} & \text{if } a \leq i \leq b \text{ and/or } a \leq j \leq b, \\ 0 & \text{otherwise,} \end{cases} \quad (11.30)$$

can be used; the rank in this case is $r = b - a + 1$, which is just the count of ones along the matrix’s main diagonal.

The name “rank- r ” implies that I_r has a “rank” of r , and indeed it does. For the moment, however, we shall discern the attribute of rank only in the rank- r identity matrix itself. Section 12.5 defines *rank* for matrices more generally.

11.3.6 The truncation operator

The rank- r identity matrix I_r is also the *truncation operator*. Attacking from the left, as in $I_r A$, it retains the first through r th rows of A but cancels other rows. Attacking from the right, as in $A I_r$, it retains the first through r th columns. Such truncation is useful symbolically to reduce an extended operator to dimension-limited form.

Whether a matrix C has dimension-limited or extended-operational form (though not necessarily if it has some other form), if it has an $m \times n$ active region¹⁷ and

$$\begin{aligned} m &\leq r, \\ n &\leq r, \end{aligned}$$

then

$$I_r C = I_r C I_r = C I_r. \quad (11.31)$$

For such a matrix, (11.31) says at least two things:

- It is superfluous to truncate both rows and columns; it suffices to truncate one or the other.
- The rank- r identity matrix I_r commutes freely past C .

Evidently big identity matrices commute freely where small ones cannot (and the general identity matrix $I = I_\infty$ commutes freely past everything).

11.3.7 The elementary vector and the lone-element matrix

The *lone-element matrix* E_{mn} is the matrix with a one in the mn th cell and zeros elsewhere:

$$[E_{mn}]_{ij} \equiv \delta_{im} \delta_{jn} = \begin{cases} 1 & \text{if } i = m \text{ and } j = n, \\ 0 & \text{otherwise.} \end{cases} \quad (11.32)$$

By this definition, $C = \sum_{i,j} c_{ij} E_{ij}$ for any matrix C . The vector analog of the lone-element matrix is the *elementary vector* \mathbf{e}_m , which has a one as the m th element:

$$[\mathbf{e}_m]_i \equiv \delta_{im} = \begin{cases} 1 & \text{if } i = m, \\ 0 & \text{otherwise.} \end{cases} \quad (11.33)$$

By this definition, $[I]_{*j} = \mathbf{e}_j$ and $[I]_{i*} = \mathbf{e}_i^T$.

¹⁷Refer to the definitions of *active region* in §§ 11.3.1 and 11.3.2. That a matrix has an $m \times n$ active region does not necessarily mean that it is all zero outside the $m \times n$ rectangle. (After all, if it were always all zero outside, then there would be little point in applying a truncation operator. There would be nothing there to truncate.)

11.3.8 Off-diagonal entries

It is interesting to observe and useful to note that if

$$[C_1]_{i*} = [C_2]_{i*} = \mathbf{e}_i^T,$$

then also

$$[C_1 C_2]_{i*} = \mathbf{e}_i^T; \quad (11.34)$$

and likewise that if

$$[C_1]_{*j} = [C_2]_{*j} = \mathbf{e}_j,$$

then also

$$[C_1 C_2]_{*j} = \mathbf{e}_j. \quad (11.35)$$

The product of matrices has off-diagonal entries in a row or column only if at least one of the factors itself has off-diagonal entries in that row or column. Or, less readably but more precisely, *the i th row or j th column of the product of matrices can depart from \mathbf{e}_i^T or \mathbf{e}_j , respectively, only if the corresponding row or column of at least one of the factors so departs*. The reason is that in (11.34), C_1 acts as a row operator on C_2 ; that if C_1 's i th row is \mathbf{e}_i^T , then its action is merely to duplicate C_2 's i th row, which itself is just \mathbf{e}_i^T . Parallel logic naturally applies to (11.35).

11.4 The elementary operator

Section 11.1.3 has introduced the general row or column operator. Conventionally denoted T , the *elementary operator* is a simple extended row or column operator from sequences of which more complicated extended operators can be built. The elementary operator T comes in three kinds.¹⁸

- The first is the *interchange elementary*

$$T_{[i \leftrightarrow j]} = I - (E_{ii} + E_{jj}) + (E_{ij} + E_{ji}), \quad (11.36)$$

which by operating $T_{[i \leftrightarrow j]}A$ or $AT_{[i \leftrightarrow j]}$ respectively interchanges A 's i th row or column with its j th.¹⁹

¹⁸In § 11.3, the symbol A specifically represented an extended operator, but here and generally the symbol represents any matrix.

¹⁹As a matter of definition, some authors [27] forbid $T_{[i \leftrightarrow i]}$ as an elementary operator, where $j = i$, since after all $T_{[i \leftrightarrow i]} = I$; which is to say that the operator doesn't actually do anything. There exist legitimate tactical reasons to forbid (as in § 11.6), but normally this book permits.

- The second is the *scaling elementary*

$$T_{\alpha[i]} = I + (\alpha - 1)E_{ii}, \quad \alpha \neq 0, \quad (11.37)$$

which by operating $T_{\alpha[i]}A$ or $AT_{\alpha[i]}$ scales (multiplies) A 's i th row or column, respectively, by the factor α .

- The third and last is the *addition elementary*

$$T_{\alpha[ij]} = I + \alpha E_{ij}, \quad i \neq j, \quad (11.38)$$

which by operating $T_{\alpha[ij]}A$ adds to the i th row of A , α times the j th row; or which by operating $AT_{\alpha[ij]}$ adds to the j th column of A , α times the i th column.

Examples of the elementary operators include

$$\begin{aligned} T_{[1 \leftrightarrow 2]} &= \begin{bmatrix} \ddots & & & & & & \\ \cdots & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 1 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & \cdots \\ & & & & & & \ddots \end{bmatrix}, \\ T_{5[4]} &= \begin{bmatrix} \ddots & & & & & & \\ \cdots & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 5 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & \cdots \\ & & & & & & \ddots \end{bmatrix}, \\ T_{5[21]} &= \begin{bmatrix} \ddots & & & & & & \\ \cdots & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 5 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 1 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & \cdots \\ & & & & & & \ddots \end{bmatrix}. \end{aligned}$$

It is good to define a concept aesthetically. One should usually do so when one can; and indeed in this case one might reasonably promote either definition on aesthetic grounds. However, an applied mathematician ought not to let a mere definition entangle him. What matters is the underlying concept. Where the definition does not serve the concept well, the applied mathematician considers whether it were not worth the effort to adapt the definition accordingly.

Note that none of these, and in fact no elementary operator of any kind, differs from I in more than four elements.

11.4.1 Properties

Significantly, elementary operators as defined above are always invertible (which is to say, reversible in effect), with

$$\begin{aligned} T_{[i \leftrightarrow j]}^{-1} &= T_{[j \leftrightarrow i]} = T_{[i \leftrightarrow j]}, \\ T_{\alpha[i]}^{-1} &= T_{(1/\alpha)[i]}, \\ T_{\alpha[ij]}^{-1} &= T_{-\alpha[ij]}, \end{aligned} \quad (11.39)$$

being themselves elementary operators such that

$$T^{-1}T = I = TT^{-1} \quad (11.40)$$

in each case.²⁰ This means that any sequence of elementaries $\prod_k T_k$ can safely be undone by the reverse sequence $\prod_k T_k^{-1}$:

$$\prod_k T_k^{-1} \prod_k T_k = I = \prod_k T_k \prod_k T_k^{-1}. \quad (11.41)$$

The rank- r identity matrix I_r is no elementary operator,²¹ nor is the lone-element matrix E_{mn} ; but the general identity matrix I is indeed an elementary operator. The last can be considered a distinct, fourth kind of elementary operator if desired; but it is probably easier just to regard it as an elementary of any of the first three kinds, since $I = T_{[i \leftrightarrow i]} = T_{1[i]} = T_{0[ij]}$.

From (11.31), we have that

$$I_r T = I_r T I_r = T I_r \quad \text{if } 1 \leq i \leq r \text{ and } 1 \leq j \leq r \quad (11.42)$$

for any elementary operator T which operates within the given bounds. Equation (11.42) lets an identity matrix with sufficiently high rank pass through a sequence of elementaries as needed.

In general, the transpose of an elementary row operator is the corresponding elementary column operator. Curiously, the interchange elementary is its own transpose and adjoint:

$$T_{[i \leftrightarrow j]}^* = T_{[i \leftrightarrow j]} = T_{[i \leftrightarrow j]}^T. \quad (11.43)$$

²⁰The addition elementary $T_{\alpha[ii]}$ and the scaling elementary $T_{0[i]}$ are forbidden precisely because they are not generally invertible.

²¹If the statement seems to contradict statements of some other books, it is only a matter of definition. This book finds it convenient to define the elementary operator in infinite-dimensional, extended-operational form. The other books are not wrong; their underlying definitions just differ slightly.

11.4.2 Commutation and sorting

Elementary operators often occur in long chains like

$$A = T_{-4[32]}T_{[2\leftrightarrow 3]}T_{(1/5)[3]}T_{(1/2)[31]}T_{5[21]}T_{[1\leftrightarrow 3]},$$

with several elementaries of all kinds intermixed. Some applications demand that the elementaries be sorted and grouped by kind, as

$$A = (T_{[2\leftrightarrow 3]}T_{[1\leftrightarrow 3]}) (T_{-4[21]}T_{(1/0 \times A)[13]}T_{5[23]}) (T_{(1/5)[1]})$$

or as

$$A = (T_{-4[32]}T_{(1/0 \times A)[21]}T_{5[31]}) (T_{(1/5)[2]}) (T_{[2\leftrightarrow 3]}T_{[1\leftrightarrow 3]}),$$

among other possible orderings. Though you probably cannot tell just by looking, the three products above are different orderings of the same elementary chain; they yield the same A and thus represent exactly the same matrix operation. Interesting is that the act of reordering the elementaries has altered some of them into other elementaries of the same kind, but has changed the kind of none of them.

One sorts a chain of elementary operators by repeatedly exchanging adjacent pairs. This of course supposes that one can exchange adjacent pairs, which seems impossible since matrix multiplication is not commutative: $A_1A_2 \neq A_2A_1$. However, at the moment we are dealing in elementary operators only; and for most pairs T_1 and T_2 of elementary operators, though indeed $T_1T_2 \neq T_2T_1$, it so happens that there exists either a T'_1 such that $T_1T_2 = T_2T'_1$ or a T'_2 such that $T_1T_2 = T'_2T_1$, where T'_1 and T'_2 are elementaries of the same kinds respectively as T_1 and T_2 . The attempt sometimes fails when both T_1 and T_2 are addition elementaries, but all other pairs commute in this way. Significantly, *elementaries of different kinds always commute*. And, though commutation can alter one (never both) of the two elementaries, it changes the kind of neither.

Many qualitatively distinct pairs of elementaries exist; we shall list these exhaustively in a moment. First, however, we should like to observe a natural hierarchy among the three kinds of elementary: (i) interchange; (ii) scaling; (iii) addition.

- The interchange elementary is the strongest. Itself subject to alteration only by another interchange elementary, it can alter any elementary by commuting past. When an interchange elementary commutes past another elementary of any kind, what it alters are the other elementary's indices i and/or j (or m and/or n , or whatever symbols

happen to represent the indices in question). When two interchange elementaries commute past one another, only one of the two is altered. (Which one? Either. The mathematician chooses.) Refer to Table 11.1.

- Next in strength is the scaling elementary. Only an interchange elementary can alter it, and it in turn can alter only an addition elementary. Scaling elementaries do not alter one another during commutation. When a scaling elementary commutes past an addition elementary, what it alters is the latter's scale α (or β , or whatever symbol happens to represent the scale in question). Refer to Table 11.2.
- The addition elementary, last and weakest, is subject to alteration by either of the other two, itself having no power to alter any elementary during commutation. A pair of addition elementaries are the only pair that can altogether fail to commute—they fail when the row index of one equals the column index of the other—but when they do commute, neither alters the other. Refer to Table 11.3.

Tables 11.1, 11.2 and 11.3 list all possible pairs of elementary operators, as the reader can check. The only pairs that fail to commute are the last three of Table 11.3.

11.5 Inversion and similarity (introduction)

If Tables 11.1, 11.2 and 11.3 exhaustively describe the commutation of one elementary past another elementary, then what can one write of the commutation of an elementary past the general matrix A ? With some matrix algebra,

$$\begin{aligned} TA &= (TA)(I) = (TA)(T^{-1}T), \\ AT &= (I)(AT) = (TT^{-1})(AT), \end{aligned}$$

one can write that

$$\begin{aligned} TA &= [TAT^{-1}]T, \\ AT &= T[T^{-1}AT], \end{aligned} \tag{11.44}$$

where T^{-1} is given by (11.39). An elementary commuting rightward changes A to TAT^{-1} ; commuting leftward, to $T^{-1}AT$.

Table 11.1: Inverting, commuting, combining and expanding elementary operators: interchange. In the table, $i \neq j \neq m \neq n$; no two indices are the same. Notice that the effect an interchange elementary $T_{[m \leftrightarrow n]}$ has in passing any other elementary, even another interchange elementary, is simply to replace m by n and n by m among the indices of the other elementary.

$$\begin{aligned}
T_{[m \leftrightarrow n]} &= T_{[n \leftrightarrow m]} \\
T_{[m \leftrightarrow m]} &= I \\
IT_{[m \leftrightarrow n]} &= T_{[m \leftrightarrow n]}I \\
T_{[m \leftrightarrow n]}T_{[m \leftrightarrow n]} &= T_{[m \leftrightarrow n]}T_{[n \leftrightarrow m]} = T_{[n \leftrightarrow m]}T_{[m \leftrightarrow n]} = I \\
T_{[m \leftrightarrow n]}T_{[i \leftrightarrow n]} &= T_{[i \leftrightarrow n]}T_{[m \leftrightarrow i]} = T_{[i \leftrightarrow m]}T_{[m \leftrightarrow n]} \\
&= (T_{[i \leftrightarrow n]}T_{[m \leftrightarrow n]})^2 \\
T_{[m \leftrightarrow n]}T_{[i \leftrightarrow j]} &= T_{[i \leftrightarrow j]}T_{[m \leftrightarrow n]} \\
T_{[m \leftrightarrow n]}T_{\alpha[m]} &= T_{\alpha[n]}T_{[m \leftrightarrow n]} \\
T_{[m \leftrightarrow n]}T_{\alpha[i]} &= T_{\alpha[i]}T_{[m \leftrightarrow n]} \\
T_{[m \leftrightarrow n]}T_{\alpha[ij]} &= T_{\alpha[ij]}T_{[m \leftrightarrow n]} \\
T_{[m \leftrightarrow n]}T_{\alpha[in]} &= T_{\alpha[im]}T_{[m \leftrightarrow n]} \\
T_{[m \leftrightarrow n]}T_{\alpha[mj]} &= T_{\alpha[nj]}T_{[m \leftrightarrow n]} \\
T_{[m \leftrightarrow n]}T_{\alpha[mn]} &= T_{\alpha[nm]}T_{[m \leftrightarrow n]}
\end{aligned}$$

Table 11.2: Inverting, commuting, combining and expanding elementary operators: scaling. In the table, $i \neq j \neq m \neq n$; no two indices are the same.

$$\begin{aligned}
T_{1[m]} &= I \\
IT_{\beta[m]} &= T_{\beta[m]}I \\
T_{(1/\beta)[m]}T_{\beta[m]} &= I \\
T_{\beta[m]}T_{\alpha[m]} &= T_{\alpha[m]}T_{\beta[m]} = T_{\alpha\beta[m]} \\
T_{\beta[m]}T_{\alpha[i]} &= T_{\alpha[i]}T_{\beta[m]} \\
T_{\beta[m]}T_{\alpha[ij]} &= T_{\alpha[ij]}T_{\beta[m]} \\
T_{\beta[m]}T_{\alpha\beta[im]} &= T_{\alpha[im]}T_{\beta[m]} \\
T_{\beta[m]}T_{\alpha[mj]} &= T_{\alpha\beta[mj]}T_{\beta[m]}
\end{aligned}$$

Table 11.3: Inverting, commuting, combining and expanding elementary operators: addition. In the table, $i \neq j \neq m \neq n$; no two indices are the same. The last three lines give pairs of addition elementaries that do not commute.

$$\begin{aligned}
T_{0[ij]} &= I \\
IT_{\alpha[ij]} &= T_{\alpha[ij]}I \\
T_{-\alpha[ij]}T_{\alpha[ij]} &= I \\
T_{\beta[ij]}T_{\alpha[ij]} &= T_{\alpha[ij]}T_{\beta[ij]} = T_{(\alpha+\beta)[ij]} \\
T_{\beta[mj]}T_{\alpha[ij]} &= T_{\alpha[ij]}T_{\beta[mj]} \\
T_{\beta[in]}T_{\alpha[ij]} &= T_{\alpha[ij]}T_{\beta[in]} \\
T_{\beta[mn]}T_{\alpha[ij]} &= T_{\alpha[ij]}T_{\beta[mn]} \\
T_{\beta[mi]}T_{\alpha[ij]} &= T_{\alpha[ij]}T_{\alpha\beta[mj]}T_{\beta[mi]} \\
T_{\beta[jn]}T_{\alpha[ij]} &= T_{\alpha[ij]}T_{-\alpha\beta[in]}T_{\beta[jn]} \\
T_{\beta[ji]}T_{\alpha[ij]} &\neq T_{\alpha[ij]}T_{\beta[ji]}
\end{aligned}$$

First encountered in § 11.4, the notation T^{-1} means the *inverse* of the elementary operator T , such that

$$T^{-1}T = I = TT^{-1}.$$

Matrix inversion is not for elementary operators only, though. Many more general matrices C also have inverses such that

$$C^{-1}C = I = CC^{-1}. \quad (11.45)$$

(Do all matrices have such inverses? No. For example, the null matrix has no such inverse.) The broad question of how to invert a general matrix C , we leave for Chs. 12 and 13 to address. For the moment however we should like to observe three simple rules involving matrix inversion.

First, nothing in the logic leading to (11.44) actually requires the matrix T there to be an elementary operator. Any matrix C for which C^{-1} is known can fill the role. Hence,

$$\begin{aligned} CA &= [CAC^{-1}]C, \\ AC &= C[C^{-1}AC]. \end{aligned} \quad (11.46)$$

The transformation CAC^{-1} or $C^{-1}AC$ is called a *similarity transformation*. Sections 12.2 and 14.9 speak further of this.

Second,

$$\begin{aligned} (C^T)^{-1} &= C^{-T} = (C^{-1})^T, \\ (C^*)^{-1} &= C^{-*} = (C^{-1})^*, \end{aligned} \quad (11.47)$$

where C^{-*} is condensed notation for conjugate transposition and inversion in either order and C^{-T} is of like style. Equation (11.47) is a consequence of (11.14), since for conjugate transposition

$$(C^{-1})^* C^* = [CC^{-1}]^* = [I]^* = I = [I]^* = [C^{-1}C]^* = C^* (C^{-1})^*$$

and similarly for nonconjugate transposition.

Third,

$$\left(\prod_k C_k \right)^{-1} = \prod_k C_k^{-1}. \quad (11.48)$$

This rule emerges upon repeated application of (11.45), which yields

$$\prod_k C_k^{-1} \prod_k C_k = I = \prod_k C_k \prod_k C_k^{-1}.$$

Table 11.4: Matrix inversion properties. (The properties work equally for $C^{-1(r)}$ as for C^{-1} if A honors an $r \times r$ active region. The full notation $C^{-1(r)}$ for the rank- r inverse incidentally is not standard, usually is not needed, and normally is not used.)

$$\begin{aligned} C^{-1}C &= I = CC^{-1} \\ C^{-1(r)}C &= I_r = CC^{-1(r)} \\ (C^T)^{-1} &= C^{-T} = (C^{-1})^T \\ (C^*)^{-1} &= C^{-*} = (C^{-1})^* \end{aligned}$$

$$\begin{aligned} CA &= [CAC^{-1}]C \\ AC &= C[C^{-1}AC] \\ \left(\prod_k C_k\right)^{-1} &= \prod_k C_k^{-1} \end{aligned}$$

A more limited form of the inverse exists than the infinite-dimensional form of (11.45). This is the rank- r inverse, a matrix $C^{-1(r)}$ such that

$$C^{-1(r)}C = I_r = CC^{-1(r)}. \quad (11.49)$$

The full notation $C^{-1(r)}$ is not standard and usually is not needed, since the context usually implies the rank. When so, one can abbreviate the notation to C^{-1} . In either notation, (11.47) and (11.48) apply equally for the rank- r inverse as for the infinite-dimensional inverse. Because of (11.31), (11.46) too applies for the rank- r inverse if A 's active region is limited to $r \times r$. (Section 13.2 uses the rank- r inverse to solve an exactly determined linear system. This is a famous way to use the inverse, with which many or most readers will already be familiar; but before using it so in Ch. 13, we shall first learn how to compute it reliably in Ch. 12.)

Table 11.4 summarizes.

11.6 Parity

Consider the sequence of integers or other objects $1, 2, 3, \dots, n$. By successively interchanging pairs of the objects (any pairs, not just adjacent pairs), one can achieve any desired permutation (§ 4.2.1). For example, beginning

with 1, 2, 3, 4, 5, one can achieve the permutation 3, 5, 1, 4, 2 by interchanging first the 1 and 3, then the 2 and 5.

Now contemplate all possible pairs:

$$\begin{array}{ccccccc} (1, 2) & (1, 3) & (1, 4) & \cdots & (1, n); \\ & (2, 3) & (2, 4) & \cdots & (2, n); \\ & & (3, 4) & \cdots & (3, n); \\ & & & \ddots & \vdots \\ & & & & (n-1, n). \end{array}$$

In a given permutation (like 3, 5, 1, 4, 2), some pairs will appear in correct order with respect to one another, while others will appear in incorrect order. (In 3, 5, 1, 4, 2, the pair [1, 2] appears in correct order in that the larger 2 stands to the right of the smaller 1; but the pair [1, 3] appears in incorrect order in that the larger 3 stands to the *left* of the smaller 1.) If p is the number of pairs which appear in incorrect order (in the example, $p = 6$), and if p is even, then we say that the permutation has *even* or *positive parity*; if odd, then *odd* or *negative parity*.²²

Now consider: every interchange of adjacent elements must either increment or decrement p by one, reversing parity. Why? Well, think about it. If two elements are adjacent and their order is correct, then interchanging falsifies the order, but only of that pair (no other element interposes, thus the interchange affects the ordering of no other pair). Complementarily, if the order is incorrect, then interchanging rectifies the order. Either way, an adjacent interchange alters p by exactly ± 1 , thus reversing parity.

What about nonadjacent elements? Does interchanging a pair of these reverse parity, too? To answer the question, let u and v represent the two elements interchanged, with a_1, a_2, \dots, a_m the elements lying between. Before the interchange:

$$\dots, u, a_1, a_2, \dots, a_{m-1}, a_m, v, \dots$$

After the interchange:

$$\dots, v, a_1, a_2, \dots, a_{m-1}, a_m, u, \dots$$

The interchange reverses with respect to one another just the pairs

$$\begin{array}{cccccc} (u, a_1) & (u, a_2) & \cdots & (u, a_{m-1}) & (u, a_m) \\ (a_1, v) & (a_2, v) & \cdots & (a_{m-1}, v) & (a_m, v) \\ (u, v) & & & & \end{array}$$

²²For readers who learned arithmetic in another language than English, the *even* integers are $\dots, -4, -2, 0, 2, 4, 6, \dots$; the *odd* integers are $\dots, -3, -1, 1, 3, 5, 7, \dots$

The number of pairs reversed is odd. Since each reversal alters p by ± 1 , the net change in p apparently also is odd, reversing parity. It seems that regardless of how distant the pair, *interchanging any pair of elements reverses the permutation's parity*.

The sole exception arises when an element is interchanged with itself. This does not change parity, but it does not change anything else, either, so in parity calculations we ignore it.²³ All other interchanges reverse parity.

We discuss parity in this, a chapter on matrices, because parity concerns the elementary interchange operator of § 11.4. The rows or columns of a matrix can be considered elements in a sequence. If so, then the interchange operator $T_{[i \leftrightarrow j]}$, $i \neq j$, acts precisely in the manner described, interchanging rows or columns and thus reversing parity. It follows that if $i_k \neq j_k$ and q is odd, then $\prod_{k=1}^q T_{[i_k \leftrightarrow j_k]} \neq I$. However, it is possible that $\prod_{k=1}^q T_{[i_k \leftrightarrow j_k]} = I$ if q is even. In any event, even q implies even p , which means even (positive) parity; odd q implies odd p , which means odd (negative) parity.

We shall have more to say about parity in §§ 11.7.1 and 14.1.

11.7 The quasidelementary operator

Multiplying sequences of the elementary operators of § 11.4, one can form much more complicated operators, which per (11.41) are always invertible. Such complicated operators are not trivial to analyze, however, so one finds it convenient to define an intermediate class of operators, called in this book the *quasidelementary operators*, more complicated than elementary operators but less so than arbitrary matrices.

A quasidelementary operator is composed of elementaries only of a single kind. There are thus three kinds of quasidelementary—interchange, scaling and addition—to match the three kinds of elementary. With respect to interchange and scaling, any sequences of elementaries of the respective kinds are allowed. With respect to addition, there are some extra rules, explained in § 11.7.3.

The three subsections which follow respectively introduce the three kinds of quasidelementary operator.

²³This is why some authors forbid self-interchanges, as explained in footnote 19.

11.7.1 The interchange quasidelementary or general interchange operator

Any product P of zero or more interchange elementaries,

$$P = \prod_k T_{[i_k \leftrightarrow j_k]}, \quad (11.50)$$

constitutes an *interchange quasidelementary, permutation matrix, permutor or general interchange operator*.²⁴ An example is

$$P = T_{[2 \leftrightarrow 5]} T_{[1 \leftrightarrow 3]} = \begin{bmatrix} \cdot & \cdot & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & \cdots \\ \cdots & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 1 & 0 & \cdots \\ \cdots & 0 & 1 & 0 & 0 & 0 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}.$$

This operator resembles I in that it has a single one in each row and in each column, but the ones here do not necessarily run along the main diagonal. The effect of the operator is to shuffle the rows or columns of the matrix it operates on, without altering any of the rows or columns it shuffles.

By (11.41), (11.39), (11.43) and (11.15), the inverse of the general interchange operator is

$$\begin{aligned} P^{-1} &= \left(\prod_k T_{[i_k \leftrightarrow j_k]} \right)^{-1} = \prod_k T_{[i_k \leftrightarrow j_k]}^{-1} \\ &= \prod_k T_{[i_k \leftrightarrow j_k]} \\ &= \prod_k T_{[i_k \leftrightarrow j_k]}^* = \left(\prod_k T_{[i_k \leftrightarrow j_k]} \right)^* \\ &= P^* = P^T \end{aligned} \quad (11.51)$$

(where $P^* = P^T$ because P has only real elements). The inverse, transpose and adjoint of the general interchange operator are thus the same:

$$P^T P = P^* P = I = P P^* = P P^T. \quad (11.52)$$

²⁴The letter P here recalls the verb “to permute.”

A significant attribute of the general interchange operator P is its parity: positive or even parity if the number of interchange elementaries $T_{[i_k \leftrightarrow j_k]}$ which compose it is even; negative or odd parity if the number is odd. This works precisely as described in § 11.6. For the purpose of parity determination, only interchange elementaries $T_{[i_k \leftrightarrow j_k]}$ for which $i_k \neq j_k$ are counted; any $T_{[i \leftrightarrow i]} = I$ noninterchanges are ignored. Thus the example P above has even parity (two interchanges), as does I itself (zero interchanges), but $T_{[i \leftrightarrow j]}$ alone (one interchange) has odd parity if $i \neq j$. As we shall see in § 14.1, the positive (even) and negative (odd) parities sometimes lend actual positive and negative senses to the matrices they describe. The parity of the general interchange operator P concerns us for this reason.

Parity, incidentally, is a property of the matrix P itself, not just of the operation P represents. No interchange quasidelementary P has positive parity as a row operator but negative as a column operator. The reason is that, regardless of whether one ultimately means to use P as a row or column operator, the matrix is nonetheless composable as a definite sequence of interchange elementaries. It is the number of interchanges, not the use, which determines P 's parity.

11.7.2 The scaling quasidelementary or general scaling operator

Like the interchange quasidelementary P of § 11.7.1, the *scaling quasidelementary, diagonal matrix* or *general scaling operator* D consists of a product of zero or more elementary operators, in this case elementary scaling operators:²⁵

$$D = \prod_{i=-\infty}^{\infty} T_{\alpha_i[i]} = \prod_{i=-\infty}^{\infty} T_{\alpha_i[i]} = \sum_{i=-\infty}^{\infty} \alpha_i E_{ii} = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & * & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & * & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & * & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & * & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & * & 0 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \end{bmatrix} \quad (11.53)$$

(of course it might be that $\alpha_i = 1$, hence that $T_{\alpha_i[i]} = I$, for some, most or even all i ; however, $\alpha_i = 0$ is forbidden by the definition of the scaling

²⁵The letter D here recalls the adjective “diagonal.”

elementary). An example is

$$D = T_{-5[4]}T_{4[2]}T_{7[1]} = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 7 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 4 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & -5 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}.$$

This operator resembles I in that all its entries run down the main diagonal; but these entries, though never zeros, are not necessarily ones, either. They are nonzero scaling factors. The effect of the operator is to scale the rows or columns of the matrix it operates on.

The general scaling operator is a particularly simple matrix. Its inverse is evidently

$$D^{-1} = \prod_{i=-\infty}^{\infty} T_{(1/\alpha_i)[i]} = \prod_{i=-\infty}^{\infty} T_{(1/\alpha_i)[i]} = \sum_{i=-\infty}^{\infty} \frac{E_{ii}}{\alpha_i}, \quad (11.54)$$

where each element down the main diagonal is individually inverted.

A superset of the general scaling operator is the *diagonal matrix*, defined less restrictively that $[A]_{ij} = 0$ for $i \neq j$, where zeros along the main diagonal are allowed. The conventional notation

$$[\text{diag}\{\mathbf{x}\}]_{ij} \equiv \delta_{ij}x_i = \delta_{ij}x_j, \quad (11.55)$$

$$\text{diag}\{\mathbf{x}\} = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & x_1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & x_2 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & x_3 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & x_4 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & x_5 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

converts a vector \mathbf{x} into a diagonal matrix. The diagonal matrix in general is not invertible and is no quasidelementary operator, but is sometimes useful nevertheless.

11.7.3 Addition quasidelementaries

Any product of interchange elementaries (§ 11.7.1), any product of scaling elementaries (§ 11.7.2), qualifies as a quasidelementary operator. Not so, any

product of addition elementaries. To qualify as a quasielementary, a product of elementary addition operators must meet some additional restrictions.

Four types of addition quasielementary are defined:²⁶

- the *downward multitarget row addition operator*,²⁷

$$\begin{aligned}
 L_{[j]} &= \prod_{i=j+1}^{\infty} T_{\alpha_{ij}[ij]} = \prod_{i=j+1}^{\infty} T_{\alpha_{ij}[ij]} \\
 &= I + \sum_{i=j+1}^{\infty} \alpha_{ij} E_{ij} \\
 &= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & * & 1 & 0 & \cdots \\ \cdots & 0 & 0 & * & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},
 \end{aligned} \tag{11.56}$$

whose inverse is

$$\begin{aligned}
 L_{[j]}^{-1} &= \prod_{i=j+1}^{\infty} T_{-\alpha_{ij}[ij]} = \prod_{i=j+1}^{\infty} T_{-\alpha_{ij}[ij]} \\
 &= I - \sum_{i=j+1}^{\infty} \alpha_{ij} E_{ij} = 2I - L_{[j]};
 \end{aligned} \tag{11.57}$$

²⁶In this subsection the explanations are briefer than in the last two, but the pattern is similar. The reader can fill in the details.

²⁷The letter L here recalls the adjective “lower.”

- the *upward multitarget row addition operator*,²⁸

$$\begin{aligned}
 U_{[j]} &= \prod_{i=-\infty}^{j-1} T_{\alpha_{ij}[ij]} = \prod_{i=-\infty}^{j-1} T_{\alpha_{ij}[ij]} \quad (11.58) \\
 &= I + \sum_{i=-\infty}^{j-1} \alpha_{ij} E_{ij} \\
 &= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 1 & 0 & * & 0 & 0 & \cdots \\ \cdots & 0 & 1 & * & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 1 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},
 \end{aligned}$$

whose inverse is

$$\begin{aligned}
 U_{[j]}^{-1} &= \prod_{i=-\infty}^{j-1} T_{-\alpha_{ij}[ij]} = \prod_{i=-\infty}^{j-1} T_{-\alpha_{ij}[ij]} \quad (11.59) \\
 &= I - \sum_{i=-\infty}^{j-1} \alpha_{ij} E_{ij} = 2I - U_{[j]};
 \end{aligned}$$

- the *rightward multitarget column addition operator*, which is the transpose $L_{[j]}^T$ of the downward operator; and
- the *leftward multitarget column addition operator*, which is the transpose $U_{[j]}^T$ of the upward operator.

11.8 The unit triangular matrix

Yet more complicated than the quasidelementary of § 11.7 is the *unit triangular matrix*, with which we draw this necessary but tedious chapter toward

²⁸The letter U here recalls the adjective “upper.”

a long close:

$$L = I + \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{i-1} \alpha_{ij} E_{ij} = I + \sum_{j=-\infty}^{\infty} \sum_{i=j+1}^{\infty} \alpha_{ij} E_{ij} \quad (11.60)$$

$$= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & * & 1 & 0 & 0 & 0 & \cdots \\ \cdots & * & * & 1 & 0 & 0 & \cdots \\ \cdots & * & * & * & 1 & 0 & \cdots \\ \cdots & * & * & * & * & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix};$$

$$U = I + \sum_{i=-\infty}^{\infty} \sum_{j=i+1}^{\infty} \alpha_{ij} E_{ij} = I + \sum_{j=-\infty}^{\infty} \sum_{i=-\infty}^{j-1} \alpha_{ij} E_{ij} \quad (11.61)$$

$$= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 1 & * & * & * & * & \cdots \\ \cdots & 0 & 1 & * & * & * & \cdots \\ \cdots & 0 & 0 & 1 & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 1 & * & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}.$$

The former is a *unit lower triangular matrix*; the latter, a *unit upper triangular matrix*. The unit triangular matrix is a generalized addition quasi-elementary, which adds not only to multiple targets but also from multiple sources—but in one direction only: downward or leftward for L or U^T (or U^*); upward or rightward for U or L^T (or L^*).

The general *triangular matrix* L_S or U_S , which by definition can have any values along its main diagonal, is sometimes of interest, as in the Schur decomposition of § 14.10.²⁹ The *strictly triangular matrix* $L - I$ or $U - I$ is likewise sometimes of interest, as in Table 11.5.³⁰ However, such matrices cannot in general be expressed as products of elementary operators and this section does not treat them.

This section presents and derives the basic properties of the unit triangular matrix.

²⁹The subscript S here stands for Schur. Other books typically use the symbols L and U for the general triangular matrix of Schur, but this book distinguishes by the subscript.

³⁰[46, “Schur decomposition,” 00:32, 30 August 2007]

11.8.1 Construction

To make a unit triangular matrix is straightforward:

$$\begin{aligned} L &= \prod_{j=-\infty}^{\infty} L_{[j]}; \\ U &= \prod_{j=-\infty}^{\infty} U_{[j]}. \end{aligned} \tag{11.62}$$

So long as the multiplication is done in the order indicated,³¹ then conveniently,

$$\begin{aligned} [L]_{ij} &= [L_{[j]}]_{ij}, \\ [U]_{ij} &= [U_{[j]}]_{ij}, \end{aligned} \tag{11.63}$$

which is to say that the entries of L and U are respectively nothing more than the relevant entries of the several $L_{[j]}$ and $U_{[j]}$. Equation (11.63) enables one to use (11.62) immediately and directly, without calculation, to build any unit triangular matrix desired.

The correctness of (11.63) is most easily seen if the several $L_{[j]}$ and $U_{[j]}$ are regarded as column operators acting sequentially on I :

$$\begin{aligned} L &= (I) \left(\prod_{j=-\infty}^{\infty} L_{[j]} \right); \\ U &= (I) \left(\prod_{j=-\infty}^{\infty} U_{[j]} \right). \end{aligned}$$

The reader can construct an inductive proof symbolically on this basis without too much difficulty if desired, but just thinking about how $L_{[j]}$ adds columns leftward and $U_{[j]}$, rightward, then considering the order in which the several $L_{[j]}$ and $U_{[j]}$ act, (11.63) follows at once.

³¹Recall again from § 2.3 that $\prod_k A_k = \cdots A_3 A_2 A_1$, whereas $\coprod_k A_k = A_1 A_2 A_3 \cdots$. This means that $(\prod_k A_k)(C)$ applies first A_1 , then A_2 , A_3 and so on, as row operators to C ; whereas $(C)(\coprod_k A_k)$ applies first A_1 , then A_2 , A_3 and so on, as column operators to C . The symbols \prod and \coprod as this book uses them can thus be thought of respectively as row and column sequencers.

11.8.2 The product of like unit triangular matrices

The product of like unit triangular matrices,

$$\begin{aligned} L_1 L_2 &= L, \\ U_1 U_2 &= U, \end{aligned} \tag{11.64}$$

is another unit triangular matrix of the same type. The proof for unit lower and unit upper triangular matrices is the same. In the unit lower triangular case, one starts from a form of the definition of a unit lower triangular matrix:

$$[L_1]_{ij} \text{ or } [L_2]_{ij} = \begin{cases} 0 & \text{if } i < j, \\ 1 & \text{if } i = j. \end{cases}$$

Then,

$$[L_1 L_2]_{ij} = \sum_{m=-\infty}^{\infty} [L_1]_{im} [L_2]_{mj}.$$

But as we have just observed, $[L_1]_{im}$ is null when $i < m$, and $[L_2]_{mj}$ is null when $m < j$. Therefore,

$$[L_1 L_2]_{ij} = \begin{cases} 0 & \text{if } i < j, \\ \sum_{m=j}^i [L_1]_{im} [L_2]_{mj} & \text{if } i \geq j. \end{cases}$$

Inasmuch as this is true, nothing prevents us from weakening the statement to read

$$[L_1 L_2]_{ij} = \begin{cases} 0 & \text{if } i < j, \\ \sum_{m=j}^i [L_1]_{im} [L_2]_{mj} & \text{if } i = j. \end{cases}$$

But this is just

$$[L_1 L_2]_{ij} = \begin{cases} 0 & \text{if } i < j, \\ [L_1]_{ij} [L_2]_{ij} = [L_1]_{ii} [L_2]_{ii} = (1)(1) = 1 & \text{if } i = j, \end{cases}$$

which again is the very definition of a unit lower triangular matrix. Hence (11.64).

11.8.3 Inversion

Inasmuch as any unit triangular matrix can be constructed from addition quasiaelementaries by (11.62), inasmuch as (11.63) supplies the specific quasiaelementaries, and inasmuch as (11.57) or (11.59) gives the inverse of each

such quasidelementary, one can always invert a unit triangular matrix easily by

$$\begin{aligned} L^{-1} &= \prod_{j=-\infty}^{\infty} L_{[j]}^{-1}, \\ U^{-1} &= \prod_{j=-\infty}^{\infty} U_{[j]}^{-1}. \end{aligned} \tag{11.65}$$

In view of (11.64), therefore, *the inverse of a unit lower triangular matrix is another unit lower triangular matrix; and the inverse of a unit upper triangular matrix, another unit upper triangular matrix.*

It is plain to see but still interesting to note that—unlike the inverse—the adjoint or transpose of a unit lower triangular matrix is a unit upper triangular matrix; and that the adjoint or transpose of a unit upper triangular matrix is a unit lower triangular matrix. The adjoint reverses the sense of the triangle.

11.8.4 The parallel unit triangular matrix

If a unit triangular matrix fits the special, restricted form

$$\begin{aligned} L_{\parallel}^{\{k\}} &= I + \sum_{j=-\infty}^k \sum_{i=k+1}^{\infty} \alpha_{ij} E_{ij} \\ &= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & * & * & * & 1 & 0 & \cdots \\ \cdots & * & * & * & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix} \end{aligned} \tag{11.66}$$

or

$$\begin{aligned}
 U_{\parallel}^{\{k\}} &= I + \sum_{j=k}^{\infty} \sum_{i=-\infty}^{k-1} \alpha_{ij} E_{ij} \\
 &= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 1 & 0 & * & * & * & \cdots \\ \cdots & 0 & 1 & * & * & * & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 1 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},
 \end{aligned} \tag{11.67}$$

confining its nonzero elements to a rectangle within the triangle as shown, then it is a *parallel unit triangular matrix* and has some special properties the general unit triangular matrix lacks.

The general unit lower triangular matrix L acting LA on a matrix A adds the rows of A downward. The parallel unit lower triangular matrix $L_{\parallel}^{\{k\}}$ acting $L_{\parallel}^{\{k\}}A$ also adds rows downward, but with the useful restriction that it makes no row of A both source and target. The addition is *from* A 's rows through the k th, *to* A 's $(k+1)$ th row onward. A horizontal frontier separates source from target, which thus march in A as separate squads.

Similar observations naturally apply with respect to the parallel unit upper triangular matrix $U_{\parallel}^{\{k\}}$, which acting $U_{\parallel}^{\{k\}}A$ adds rows upward, and also with respect to $L_{\parallel}^{\{k\}T}$ and $U_{\parallel}^{\{k\}T}$, which acting $AL_{\parallel}^{\{k\}T}$ and $AU_{\parallel}^{\{k\}T}$ add columns respectively rightward and leftward (remembering that $L_{\parallel}^{\{k\}T}$ is no unit lower but a unit upper triangular matrix; that $U_{\parallel}^{\{k\}T}$ is the lower). Each separates source from target in the matrix A it operates on.

The reason we care about the separation of source from target is that, in matrix arithmetic generally, where source and target are not separate but remain intermixed, the sequence matters in which rows or columns are added. That is, in general,

$$T_{\alpha_1[i_1j_1]}T_{\alpha_2[i_2j_2]} \neq I + \alpha_1 E_{i_1j_1} + \alpha_2 E_{i_2j_2} \neq T_{\alpha_2[i_2j_2]}T_{\alpha_1[i_1j_1]}.$$

It makes a difference whether the one addition comes before, during or after the other—but only because the target of the one addition might be the source of the other. The danger is that $i_1 = j_2$ or $i_2 = j_1$. Remove this danger, and the sequence ceases to matter (refer to Table 11.3).

That is exactly what the parallel unit triangular matrix does: it separates source from target and thus removes the danger. It is for this reason that

the parallel unit triangular matrix brings the useful property that

$$\begin{aligned}
L_{\parallel}^{\{k\}} &= I + \sum_{j=-\infty}^k \sum_{i=k+1}^{\infty} \alpha_{ij} E_{ij} \\
&= \prod_{j=-\infty}^k \prod_{i=k+1}^{\infty} T_{\alpha_{ij}[ij]} = \prod_{j=-\infty}^k \prod_{i=k+1}^{\infty} T_{\alpha_{ij}[ij]} \\
&= \prod_{j=-\infty}^k \prod_{i=k+1}^{\infty} T_{\alpha_{ij}[ij]} = \prod_{j=-\infty}^k \prod_{i=k+1}^{\infty} T_{\alpha_{ij}[ij]} \\
&= \prod_{i=k+1}^{\infty} \prod_{j=-\infty}^k T_{\alpha_{ij}[ij]} = \prod_{i=k+1}^{\infty} \prod_{j=-\infty}^k T_{\alpha_{ij}[ij]} \quad (11.68) \\
&= \prod_{i=k+1}^{\infty} \prod_{j=-\infty}^k T_{\alpha_{ij}[ij]} = \prod_{i=k+1}^{\infty} \prod_{j=-\infty}^k T_{\alpha_{ij}[ij]}, \\
U_{\parallel}^{\{k\}} &= I + \sum_{j=k}^{\infty} \sum_{i=-\infty}^{k-1} \alpha_{ij} E_{ij} \\
&= \prod_{j=k}^{\infty} \prod_{i=-\infty}^{k-1} T_{\alpha_{ij}[ij]} = \cdots,
\end{aligned}$$

which says that one can build a parallel unit triangular matrix equally well in any sequence—in contrast to the case of the general unit triangular matrix, whose construction per (11.62) one must sequence carefully. (Though eqn. 11.68 does not show them, even more sequences are possible. You can scramble the factors' ordering any random way you like. The multiplication is fully commutative.) Under such conditions, the inverse of the parallel unit

triangular matrix is particularly simple:³²

$$\begin{aligned}
 L_{\parallel}^{\{k\}-1} &= I - \sum_{j=-\infty}^k \sum_{i=k+1}^{\infty} \alpha_{ij} E_{ij} = 2I - L_{\parallel}^{\{k\}} \\
 &= \prod_{j=-\infty}^k \prod_{i=k+1}^{\infty} T_{-\alpha_{ij}[ij]} = \cdots, \\
 U_{\parallel}^{\{k\}-1} &= I - \sum_{j=k}^{\infty} \sum_{i=-\infty}^{k-1} \alpha_{ij} E_{ij} = 2I - U_{\parallel}^{\{k\}} \\
 &= \prod_{j=k}^{\infty} \prod_{i=-\infty}^{k-1} T_{-\alpha_{ij}[ij]} = \cdots,
 \end{aligned} \tag{11.69}$$

where again the elementaries can be multiplied in any order. Pictorially,

$$\begin{aligned}
 L_{\parallel}^{\{k\}-1} &= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & -* & -* & -* & 1 & 0 & \cdots \\ \cdots & -* & -* & -* & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \\
 U_{\parallel}^{\{k\}-1} &= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 1 & 0 & -* & -* & -* & \cdots \\ \cdots & 0 & 1 & -* & -* & -* & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 1 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}.
 \end{aligned}$$

The inverse of a parallel unit triangular matrix is just the matrix itself, only with each element off the main diagonal negated. Table 11.5 records a few properties that come immediately of the last observation and from the parallel unit triangular matrix's basic layout.

³²There is some odd parochiality at play in applied mathematics when one calls such collections of symbols as (11.69) “particularly simple.” Nevertheless, in the present context the idea (11.69) represents is indeed simple: that one can multiply constituent elementaries in any order and still reach the same parallel unit triangular matrix; that the elementaries in this case do not interfere.

Table 11.5: Properties of the parallel unit triangular matrix. (In the table, the notation I_a^b represents the generalized dimension-limited identity matrix or truncator of eqn. 11.30. Note that the inverses $L_{\parallel}^{\{k\}-1} = L_{\parallel}^{\{k\}'}$ and $U_{\parallel}^{\{k\}-1} = U_{\parallel}^{\{k\}'}$ are parallel unit triangular matrices themselves, such that the table's properties hold for them, too.)

$$\frac{L_{\parallel}^{\{k\}} + L_{\parallel}^{\{k\}-1}}{2} = I = \frac{U_{\parallel}^{\{k\}} + U_{\parallel}^{\{k\}-1}}{2}$$

$$\begin{aligned} I_{k+1}^{\infty} L_{\parallel}^{\{k\}} I_{-\infty}^k &= L_{\parallel}^{\{k\}} - I = I_{k+1}^{\infty} (L_{\parallel}^{\{k\}} - I) I_{-\infty}^k \\ I_{-\infty}^{k-1} U_{\parallel}^{\{k\}} I_k^{\infty} &= U_{\parallel}^{\{k\}} - I = I_{-\infty}^{k-1} (U_{\parallel}^{\{k\}} - I) I_k^{\infty} \end{aligned}$$

If $L_{\parallel}^{\{k\}}$ honors an $n \times n$ active region, then

$$\begin{aligned} (I_n - I_k) L_{\parallel}^{\{k\}} I_k &= L_{\parallel}^{\{k\}} - I = (I_n - I_k) (L_{\parallel}^{\{k\}} - I) I_k \\ \text{and } (I - I_n) (L_{\parallel}^{\{k\}} - I) &= 0 = (L_{\parallel}^{\{k\}} - I) (I - I_n). \end{aligned}$$

If $U_{\parallel}^{\{k\}}$ honors an $n \times n$ active region, then

$$\begin{aligned} I_{k-1} U_{\parallel}^{\{k\}} (I_n - I_{k-1}) &= U_{\parallel}^{\{k\}} - I = I_{k-1} (U_{\parallel}^{\{k\}} - I) (I_n - I_{k-1}) \\ \text{and } (I - I_n) (U_{\parallel}^{\{k\}} - I) &= 0 = (U_{\parallel}^{\{k\}} - I) (I - I_n). \end{aligned}$$

11.8.5 The partial unit triangular matrix

Besides the notation L and U for the general unit lower and unit upper triangular matrices and the notation $L_{\parallel}^{\{k\}}$ and $U_{\parallel}^{\{k\}}$ for the parallel unit lower and unit upper triangular matrices, we shall find it useful to introduce the additional notation

$$L^{[k]} = I + \sum_{j=k}^{\infty} \sum_{i=j+1}^{\infty} \alpha_{ij} E_{ij} \quad (11.70)$$

$$= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & * & 1 & 0 & \cdots \\ \cdots & 0 & 0 & * & * & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

$$U^{[k]} = I + \sum_{j=-\infty}^k \sum_{i=-\infty}^{j-1} \alpha_{ij} E_{ij} \quad (11.71)$$

$$= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & * & * & 0 & 0 & \cdots \\ \cdots & 0 & 1 & * & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 1 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

for unit triangular matrices whose off-diagonal content is confined to a narrow wedge and

$$L^{\{k\}} = I + \sum_{j=-\infty}^k \sum_{i=j+1}^{\infty} \alpha_{ij} E_{ij} \quad (11.72)$$

$$= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 1 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & * & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & * & * & 1 & 0 & 0 & 0 & \cdots \\ \cdots & * & * & * & 1 & 0 & 0 & \cdots \\ \cdots & * & * & * & 0 & 1 & 0 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

$$U^{\{k\}} = I + \sum_{j=k}^{\infty} \sum_{i=-\infty}^{j-1} \alpha_{ij} E_{ij} \quad (11.73)$$

$$= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 1 & 0 & * & * & * & * & \cdots \\ \cdots & 0 & 1 & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 1 & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 1 & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & * & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

for the supplementary forms.³³ Such notation is not standard in the literature, but it serves a purpose in this book and is introduced here for this reason. If names are needed for $L^{[k]}$, $U^{[k]}$, $L^{\{k\}}$ and $U^{\{k\}}$, the former pair can be called *minor partial unit triangular matrices*, and the latter pair, *major partial unit triangular matrices*. Whether minor or major, the partial unit triangular matrix is a matrix which leftward or rightward of the k th column resembles I . Of course partial unit triangular matrices which resemble I above or below the k th row are equally possible, and can be denoted $L^{[k]T}$, $U^{[k]T}$, $L^{\{k\}T}$ and $U^{\{k\}T}$.

Observe that the parallel unit triangular matrices $L_{\parallel}^{\{k\}}$ and $U_{\parallel}^{\{k\}}$ of § 11.8.4 are in fact also major partial unit triangular matrices, as the notation suggests.

³³The notation is arguably imperfect in that $L^{\{k\}} + L^{[k]} - I \neq L$ but rather that $L^{\{k\}} + L^{[k+1]} - I = L$. The conventional notation $\sum_{k=a}^b f(k) + \sum_{k=b}^c f(k) \neq \sum_{k=a}^c f(k)$ suffers the same arguable imperfection.

11.9 The shift operator

Not all useful matrices fit the dimension-limited and extended-operational forms of § 11.3. An exception is the *shift operator* H_k , defined that

$$[H_k]_{ij} = \delta_{i(j+k)}. \quad (11.74)$$

For example,

$$H_2 = \begin{bmatrix} & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 1 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 1 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & 0 & 0 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \end{bmatrix}.$$

Operating $H_k A$, H_k shifts A 's rows downward k steps. Operating $A H_k$, H_k shifts A 's columns leftward k steps. Inasmuch as the shift operator shifts all rows or columns of the matrix it operates on, its active region is $\infty \times \infty$ in extent. Obviously, the shift operator's inverse, transpose and adjoint are the same:

$$\begin{aligned} H_k^T H_k &= H_k^* H_k = I = H_k H_k^* = H_k H_k^T, \\ H_k^{-1} &= H_k^T = H_k^* = H_{-k}. \end{aligned} \quad (11.75)$$

Further obvious but useful identities include that

$$\begin{aligned} (I_\ell - I_k) H_k &= H_k I_{\ell-k}, \\ H_{-k} (I_\ell - I_k) &= I_{\ell-k} H_{-k}. \end{aligned} \quad (11.76)$$

11.10 The Jacobian derivative

Chapter 4 has introduced the derivative of a function with respect to a scalar variable. One can also take the derivative of a function with respect to a vector variable, and the function itself can be vector-valued. The derivative is

$$\left[\frac{d\mathbf{f}}{d\mathbf{x}} \right]_{ij} = \frac{\partial f_i}{\partial x_j}. \quad (11.77)$$

For instance, if \mathbf{x} has three elements and \mathbf{f} has two, then

$$\frac{d\mathbf{f}}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} \end{bmatrix}.$$

This is called the *Jacobian derivative*, the *Jacobian matrix*, or just the *Jacobian*.³⁴ Each of its columns is the derivative with respect to one element of \mathbf{x} .

The Jacobian derivative of a vector with respect to itself is

$$\frac{d\mathbf{x}}{d\mathbf{x}} = I. \quad (11.78)$$

The derivative is not I_n as one might think, because, even if \mathbf{x} has only n elements, still, one could vary x_{n+1} in principle, and $\partial x_{n+1}/\partial x_{n+1} \neq 0$.

The Jacobian derivative obeys the derivative product rule (4.25) in the form³⁵

$$\begin{aligned} \frac{d}{d\mathbf{x}}(\mathbf{g}^T A \mathbf{f}) &= \left[\mathbf{g}^T A \left(\frac{d\mathbf{f}}{d\mathbf{x}} \right) \right] + \left[\left(\frac{d\mathbf{g}}{d\mathbf{x}} \right)^T A \mathbf{f} \right]^T, \\ \frac{d}{d\mathbf{x}}(\mathbf{g}^* A \mathbf{f}) &= \left[\mathbf{g}^* A \left(\frac{d\mathbf{f}}{d\mathbf{x}} \right) \right] + \left[\left(\frac{d\mathbf{g}}{d\mathbf{x}} \right)^* A \mathbf{f} \right]^T, \end{aligned} \quad (11.79)$$

valid for any constant matrix A —as is seen by applying the definition (4.19) of the derivative, which here is

$$\frac{\partial(\mathbf{g}^* A \mathbf{f})}{\partial x_j} = \lim_{\partial x_j \rightarrow 0} \frac{(\mathbf{g} + \partial \mathbf{g}/2)^* A (\mathbf{f} + \partial \mathbf{f}/2) - (\mathbf{g} - \partial \mathbf{g}/2)^* A (\mathbf{f} - \partial \mathbf{f}/2)}{\partial x_j},$$

and simplifying.

The shift operator of § 11.9 and the Jacobian derivative of this section complete the family of matrix rudiments we shall need to begin to do increasingly interesting things with matrices in Chs. 13 and 14. Before doing interesting things, however, we must treat two more foundational matrix matters. The two are the Gauss-Jordan decomposition and the matter of matrix rank, which will be the subjects of Ch. 12, next.

³⁴[46, “Jacobian,” 00:50, 15 Sept. 2007]

³⁵Notice that the last term on (11.79)’s second line is transposed, not adjointed.

Chapter 12

Matrix rank and the Gauss-Jordan decomposition

Chapter 11 has brought the matrix and its rudiments, the latter including

- lone-element matrix E (§ 11.3.7),
- the null matrix 0 (§ 11.3.1),
- the rank- r identity matrix I_r (§ 11.3.5),
- the general identity matrix I and the scalar matrix λI (§ 11.3.2),
- the elementary operator T (§ 11.4),
- the quasielementary operator P , D , $L_{[k]}$ or $U_{[k]}$ (§ 11.7), and
- the unit triangular matrix L or U (§ 11.8).

Such rudimentary forms have useful properties, as we have seen. The general matrix A does not necessarily have any of these properties, but it turns out that one can factor any matrix whatsoever into a product of rudiments which do have the properties, and that several orderly procedures are known to do so. The simplest of these, and indeed one of the more useful, is the Gauss-Jordan decomposition. This chapter introduces it.

Section 11.3 has de-emphasized the concept of matrix dimensionality $m \times n$, supplying in its place the new concept of matrix rank. However, that section has actually defined rank only for the rank- r identity matrix I_r . In fact all matrices have rank. This chapter explains.

Before treating the Gauss-Jordan decomposition and the matter of matrix rank as such, however, we shall find it helpful to prepare two preliminaries thereto: (i) the matter of the linear independence of vectors; and (ii) the elementary similarity transformation. The chapter begins with these.

Except in § 12.2, the chapter demands more rigor than one likes in such a book as this. However, it is hard to see how to avoid the rigor here, and logically the chapter cannot be omitted. We shall drive through the chapter in as few pages as can be managed, and then onward to the more interesting matrix topics of Chs. 13 and 14.

12.1 Linear independence

Linear independence is a significant possible property of a set of vectors—whether the set be the several columns of a matrix, the several rows, or some other vectors—the property being defined as follows. A vector is *linearly independent* if its role cannot be served by the other vectors in the set. More formally, the n vectors of the set $\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4, \mathbf{a}_5, \dots, \mathbf{a}_n\}$ are linearly independent if and only if none of them can be expressed as a linear combination—a weighted sum—of the others. That is, the several \mathbf{a}_k are linearly independent iff

$$\alpha_1 \mathbf{a}_1 + \alpha_2 \mathbf{a}_2 + \alpha_3 \mathbf{a}_3 + \cdots + \alpha_n \mathbf{a}_n \neq 0 \quad (12.1)$$

for all nontrivial α_k , where “nontrivial α_k ” means the several α_k , at least one of which is nonzero (*trivial* α_k , by contrast, would be $\alpha_1 = \alpha_2 = \alpha_3 = \cdots = \alpha_n = 0$). Vectors which can combine nontrivially to reach the null vector are by definition *linearly dependent*.

Linear independence is a property of vectors. Technically the property applies to scalars, too, inasmuch as a scalar resembles a one-element vector—so, any nonzero scalar alone is linearly independent—but there is no such thing as a linearly independent pair of scalars, because one of the pair can always be expressed as a complex multiple of the other. Significantly but less obviously, there is also no such thing as a linearly independent set which includes the null vector; (12.1) forbids it. Paradoxically, even the single-member, $n = 1$ set consisting only of $\mathbf{a}_1 = 0$ is, strictly speaking, not linearly independent.

For consistency of definition, we regard the empty, $n = 0$ set as linearly independent, on the technical ground that the only possible linear combination of the empty set is trivial.¹

¹This is the kind of thinking which typically governs mathematical edge cases. One

If a linear combination of several independent vectors \mathbf{a}_k forms a vector \mathbf{b} , then one might ask: can there exist a different linear combination of the same vectors \mathbf{a}_k which also forms \mathbf{b} ? That is, if

$$\beta_1 \mathbf{a}_1 + \beta_2 \mathbf{a}_2 + \beta_3 \mathbf{a}_3 + \cdots + \beta_n \mathbf{a}_n = \mathbf{b},$$

where the several \mathbf{a}_k satisfy (12.1), then is

$$\beta'_1 \mathbf{a}_1 + \beta'_2 \mathbf{a}_2 + \beta'_3 \mathbf{a}_3 + \cdots + \beta'_n \mathbf{a}_n = \mathbf{b}$$

possible? To answer the question, suppose that it were possible. The difference of the two equations then would be

$$(\beta'_1 - \beta_1) \mathbf{a}_1 + (\beta'_2 - \beta_2) \mathbf{a}_2 + (\beta'_3 - \beta_3) \mathbf{a}_3 + \cdots + (\beta'_n - \beta_n) \mathbf{a}_n = \mathbf{0}.$$

According to (12.1), this could only be so if the coefficients in the last equation were trivial—that is, only if $\beta'_1 - \beta_1 = 0$, $\beta'_2 - \beta_2 = 0$, $\beta'_3 - \beta_3 = 0$, \dots , $\beta'_n - \beta_n = 0$. But this says no less than that the two linear combinations, which we had supposed to differ, were in fact one and the same. One concludes therefore that, *if a vector \mathbf{b} can be expressed as a linear combination of several linearly independent vectors \mathbf{a}_k , then it cannot be expressed as any other combination of the same vectors.* The combination is unique.

Linear independence can apply in any dimensionality, but it helps to visualize the concept geometrically in three dimensions, using the three-dimensional geometrical vectors of § 3.3. Two such vectors are independent so long as they do not lie along the same line. A third such vector is independent of the first two so long as it does not lie in their common plane. A fourth such vector (unless it points off into some unvisualizable fourth dimension) cannot possibly then be independent of the three.

We discuss the linear independence of vectors in this, a chapter on matrices, because (§ 11.1) a matrix is essentially a sequence of vectors—either of column vectors or of row vectors, depending on one's point of view. As we shall see in § 12.5, the important property of matrix *rank* depends on the number of linearly independent columns or rows a matrix has.

could define the empty set to be linearly dependent if one really wanted to, but what then of the observation that adding a vector to a linearly dependent set never renders the set independent? Surely in this light it is preferable just to define the empty set as independent in the first place. Similar thinking makes $0! = 1$, $\sum_{k=0}^{-1} a_k z^k = 0$, and 2 not 1 the least prime, among other examples.

12.2 The elementary similarity transformation

Section 11.5 and its (11.46) have introduced the *similarity transformation* CAC^{-1} or $C^{-1}AC$, which arises when an operator C commutes respectively rightward or leftward past a matrix A . The similarity transformation has several interesting properties, some of which we are now prepared to discuss, particularly in the case in which the operator happens to be an elementary, $C = T$. In this case, the several rules of Table 12.1 obtain.

Most of the table's rules are fairly obvious if the meaning of the symbols is understood, though to grasp some of the rules it helps to sketch the relevant matrices on a sheet of paper. Of course rigorous symbolic proofs can be constructed after the pattern of § 11.8.2, but they reveal little or nothing sketching the matrices does not. The symbols P , D , L and U of course represent the quasidelementaries and unit triangular matrices of §§ 11.7 and 11.8. The symbols P' , D' , L' and U' also represent quasidelementaries and unit triangular matrices, only not necessarily the same ones P , D , L and U do.

The rules permit one to commute some but not all elementaries past a quasidelementary or unit triangular matrix, without fundamentally altering the character of the quasidelementary or unit triangular matrix, and sometimes without altering the matrix at all. The rules find use among other places in the Gauss-Jordan decomposition of § 12.3.

12.3 The Gauss-Jordan decomposition

The *Gauss-Jordan decomposition* of an arbitrary, dimension-limited, $m \times n$ matrix A is²

$$\begin{aligned} A &= G_{>} I_r G_{<} \\ &= PDLU I_r KS, \\ G_{<} &\equiv KS, \\ G_{>} &\equiv PDLU, \end{aligned} \tag{12.2}$$

where

- P and S are general interchange operators (§ 11.7.1);

²Most introductory linear algebra texts this writer has met call the Gauss-Jordan decomposition instead the “ LU decomposition” and include fewer factors in it, typically merging D into L and omitting K and S . They also omit I_r , since their matrices have pre-defined dimensionality. Perhaps the reader will agree that the decomposition is cleaner as presented here.

Table 12.1: Some elementary similarity transformations.

$$\begin{aligned}
T_{[i \leftrightarrow j]} I T_{[i \leftrightarrow j]} &= I \\
T_{[i \leftrightarrow j]} P T_{[i \leftrightarrow j]} &= P' \\
T_{[i \leftrightarrow j]} D T_{[i \leftrightarrow j]} &= D' = D + ([D]_{jj} - [D]_{ii}) E_{ii} + ([D]_{ii} - [D]_{jj}) E_{jj} \\
T_{[i \leftrightarrow j]} D T_{[i \leftrightarrow j]} &= D \quad \text{if } [D]_{ii} = [D]_{jj} \\
T_{[i \leftrightarrow j]} L^{[k]} T_{[i \leftrightarrow j]} &= L^{[k]} \quad \text{if } i < k \text{ and } j < k \\
T_{[i \leftrightarrow j]} U^{[k]} T_{[i \leftrightarrow j]} &= U^{[k]} \quad \text{if } i > k \text{ and } j > k \\
T_{[i \leftrightarrow j]} L^{\{k\}'} T_{[i \leftrightarrow j]} &= L^{\{k\}'} \quad \text{if } i > k \text{ and } j > k \\
T_{[i \leftrightarrow j]} U^{\{k\}'} T_{[i \leftrightarrow j]} &= U^{\{k\}'} \quad \text{if } i < k \text{ and } j < k \\
T_{[i \leftrightarrow j]} L_{\parallel}^{\{k\}'} T_{[i \leftrightarrow j]} &= L_{\parallel}^{\{k\}'} \quad \text{if } i > k \text{ and } j > k \\
T_{[i \leftrightarrow j]} U_{\parallel}^{\{k\}'} T_{[i \leftrightarrow j]} &= U_{\parallel}^{\{k\}'} \quad \text{if } i < k \text{ and } j < k \\
T_{\alpha[i]} I T_{(1/\alpha)[i]} &= I \\
T_{\alpha[i]} D T_{(1/\alpha)[i]} &= D \\
T_{\alpha[i]} A T_{(1/\alpha)[i]} &= A' \quad \text{where } A \text{ is any of} \\
&\quad L, U, L^{[k]}, U^{[k]}, L^{\{k\}'}, U^{\{k\}'}, L_{\parallel}^{\{k\}'}, U_{\parallel}^{\{k\}'} \\
T_{\alpha[ij]} I T_{-\alpha[ij]} &= I \\
T_{\alpha[ij]} D T_{-\alpha[ij]} &= D + ([D]_{jj} - [D]_{ii}) \alpha E_{ij} \neq D' \\
T_{\alpha[ij]} D T_{-\alpha[ij]} &= D \quad \text{if } [D]_{ii} = [D]_{jj} \\
T_{\alpha[ij]} L T_{-\alpha[ij]} &= L' \quad \text{if } i > j \\
T_{\alpha[ij]} U T_{-\alpha[ij]} &= U' \quad \text{if } i < j
\end{aligned}$$

- D is a general scaling operator (§ 11.7.2);
- L and U are respectively unit lower and unit upper triangular matrices (§ 11.8);
- $K = L_{\parallel}^{\{r\}T}$ is the transpose of a parallel unit lower triangular matrix, being thus a parallel unit upper triangular matrix (§ 11.8.4);
- $G_{>}$ and $G_{<}$ are composites³ as defined by (12.2); and
- r is an unspecified rank.

The Gauss-Jordan decomposition is also called the *Gauss-Jordan factorization*.

Whether all possible matrices A have a Gauss-Jordan decomposition (they do, in fact) is a matter this section addresses. However—at least for matrices which do have one—because $G_{>}$ and $G_{<}$ are composed of invertible factors, one can left-multiply the equation $A = G_{>}I_rG_{<}$ by $G_{>}^{-1}$ and right-multiply it by $G_{<}^{-1}$ to obtain

$$\begin{aligned} U^{-1}L^{-1}D^{-1}P^{-1}AS^{-1}K^{-1} &= \\ G_{>}^{-1}AG_{<}^{-1} &= I_r, \\ S^{-1}K^{-1} &= G_{<}^{-1}, \\ U^{-1}L^{-1}D^{-1}P^{-1} &= G_{>}^{-1}, \end{aligned} \tag{12.3}$$

the Gauss-Jordan's complementary form.

12.3.1 Motive

Equation (12.2) seems inscrutable. The equation itself is easy enough to read, but just as there are many ways to factor a scalar ($0xC = [4][3] = [2]^2[3] = [2][6]$, for example), there are likewise many ways to factor a matrix. Why choose this particular way?

There are indeed many ways. We shall meet some of the others in §§ 13.10, 14.6, 14.10 and 14.12. The Gauss-Jordan decomposition we meet here however has both significant theoretical properties and useful practical applications, and in any case needs less advanced preparation to appreciate

³One can pronounce $G_{>}$ and $G_{<}$ respectively as “ G acting rightward” and “ G acting leftward.” The letter G itself can be regarded as standing for “Gauss-Jordan,” but admittedly it is chosen as much because otherwise we were running out of available Roman capitals!

than the others, and (at least as developed in this book) precedes the others logically. It emerges naturally when one posits a pair of square, $n \times n$ matrices, A and A^{-1} , for which $A^{-1}A = I_n$, where A is known and A^{-1} is to be determined. (The A^{-1} here is the $A^{-1(n)}$ of eqn. 11.49. However, it is only supposed here that $A^{-1}A = I_n$; it is not *yet* claimed that $AA^{-1} = I_n$.)

To determine A^{-1} is not an entirely trivial problem. The matrix A^{-1} such that $A^{-1}A = I_n$ may or may not exist (usually it does exist if A is square, but even then it may not, as we shall soon see), and even if it does exist, how to determine it is not immediately obvious. And still, if one can determine A^{-1} , that is only for square A ; what if A is not square? In the present subsection however we are not trying to prove anything, only to motivate, so for the moment let us suppose a square A for which A^{-1} does exist, and let us seek A^{-1} by left-multiplying A by a sequence $\prod T$ of elementary row operators, each of which makes the matrix more nearly resemble I_n . When I_n is finally achieved, then we shall have that

$$\left(\prod T\right)(A) = I_n,$$

or, left-multiplying by I_n and observing that $I_n^2 = I_n$,

$$(I_n)\left(\prod T\right)(A) = I_n,$$

which implies that

$$A^{-1} = (I_n)\left(\prod T\right).$$

The product of elementaries which transforms A to I_n , truncated (§ 11.3.6) to $n \times n$ dimensionality, itself constitutes A^{-1} . This observation is what motivates the Gauss-Jordan decomposition.

By successive steps,⁴ a concrete example:

$$\begin{aligned}
 A &= \begin{bmatrix} 2 & -4 \\ 3 & -1 \end{bmatrix}, \\
 \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix} A &= \begin{bmatrix} 1 & -2 \\ 3 & -1 \end{bmatrix}, \\
 \begin{bmatrix} 1 & 0 \\ -3 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix} A &= \begin{bmatrix} 1 & -2 \\ 0 & 5 \end{bmatrix}, \\
 \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{5} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -3 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix} A &= \begin{bmatrix} 1 & -2 \\ 0 & 1 \end{bmatrix}, \\
 \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{5} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -3 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix} A &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \\
 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{5} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -3 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix} A &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.
 \end{aligned}$$

Hence,

$$A^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{5} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -3 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} -\frac{1}{A} & \frac{2}{5} \\ -\frac{3}{A} & \frac{1}{5} \end{bmatrix}.$$

Using the elementary commutation identity that $T_{\beta[m]}T_{\alpha[mj]} = T_{\alpha\beta[mj]}T_{\beta[m]}$, from Table 11.2, to group like operators, we have that

$$A^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -\frac{3}{5} & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{5} \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} -\frac{1}{A} & \frac{2}{5} \\ -\frac{3}{A} & \frac{1}{5} \end{bmatrix};$$

or, multiplying the two scaling elementaries to merge them into a single general scaling operator (§ 11.7.2),

$$A^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -\frac{3}{5} & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{5} \end{bmatrix} = \begin{bmatrix} -\frac{1}{A} & \frac{2}{5} \\ -\frac{3}{A} & \frac{1}{5} \end{bmatrix}.$$

The last equation is written symbolically as

$$A^{-1} = I_2 U^{-1} L^{-1} D^{-1},$$

from which

$$A = D L U I_2 = \begin{bmatrix} 2 & 0 \\ 0 & 5 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \frac{3}{5} & 1 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & -4 \\ 3 & -1 \end{bmatrix}.$$

⁴Theoretically, all elementary operators including the ones here have extended-operational form (§ 11.3.2), but all those \cdots ellipses clutter the page too much. Only the 2×2 active regions are shown here.

Now, admittedly, the equation $A = DLU I_2$ is not (12.2)—or rather, it is (12.2), but only in the special case that $r = 2$ and $P = S = K = I$ —which begs the question: why do we need the factors P , S and K in the first place? The answer regarding P and S is that these factors respectively gather row and column interchange elementaries, of which the example given has used none but which other examples sometimes need or want, particularly to avoid dividing by zero when they encounter a zero in an inconvenient cell of the matrix (the reader might try reducing $A = [0 \ 1; 1 \ 0]$ to I_2 , for instance; a row or column interchange is needed here). Regarding K , this factor comes into play when A has broad rectangular rather than square shape, and also sometimes when one of the rows of A happens to be a linear combination of the others. The last point, we are not quite ready to detail yet, but at present we are only motivating not proving, so if the reader will accept the other factors and suspend judgment on K until the actual need for it emerges in § 12.3.3, step 12, then we shall proceed on this basis.

12.3.2 Method

The Gauss-Jordan decomposition of a matrix A is not discovered at one stroke but rather is gradually built up, elementary by elementary. It begins with the equation

$$A = IIIIAII,$$

where the six I hold the places of the six Gauss-Jordan factors P , D , L , U , K and S of (12.2). By successive elementary operations, the A on the right is gradually transformed into I_r , while the six I are gradually transformed into the six Gauss-Jordan factors. The decomposition thus ends with the equation

$$A = PDLUI_rKS,$$

which is (12.2). In between, while the several matrices are gradually being transformed, the equation is represented as

$$A = \tilde{P}\tilde{D}\tilde{L}\tilde{U}\tilde{I}\tilde{K}\tilde{S}, \tag{12.4}$$

where the initial value of \tilde{I} is A and the initial values of \tilde{P} , \tilde{D} , etc., are all I .

Each step of the transformation goes as follows. The matrix \tilde{I} is left- or right-multiplied by an elementary operator T . To compensate, one of the six factors is right- or left-multiplied by T^{-1} . Intervening factors are multiplied by both T and T^{-1} , which multiplication constitutes an elementary similarity transformation as described in § 12.2. For example,

$$A = \tilde{P} \left(\tilde{D} T_{(1/\alpha)[i]} \right) \left(T_{\alpha[i]} \tilde{L} T_{(1/\alpha)[i]} \right) \left(T_{\alpha[i]} \tilde{U} T_{(1/\alpha)[i]} \right) \left(T_{\alpha[i]} \tilde{I} \right) \tilde{K} \tilde{S},$$

which is just (12.4), inasmuch as the adjacent elementaries cancel one another; then,

$$\begin{aligned}\tilde{I} &\leftarrow T_{\alpha[i]}\tilde{I}, \\ \tilde{U} &\leftarrow T_{\alpha[i]}\tilde{U}T_{(1/\alpha)[i]}, \\ \tilde{L} &\leftarrow T_{\alpha[i]}\tilde{L}T_{(1/\alpha)[i]}, \\ \tilde{D} &\leftarrow \tilde{D}T_{(1/\alpha)[i]},\end{aligned}$$

thus associating the operation with the appropriate factor—in this case, \tilde{D} . Such elementary row and column operations are repeated until $\tilde{I} = I_r$, at which point (12.4) has become the Gauss-Jordan decomposition (12.2).

12.3.3 The algorithm

Having motivated the Gauss-Jordan decomposition in § 12.3.1 and having proposed a basic method to pursue it in § 12.3.2, we shall now establish a definite, orderly, failproof algorithm to achieve it. Broadly, the algorithm

- copies A into the variable working matrix \tilde{I} (step 1 below),
- reduces \tilde{I} by suitable row (and maybe column) operations to unit upper triangular form (steps 2 through 7),
- establishes a rank r (step 8), and
- reduces the now unit triangular \tilde{I} further to the rank- r identity matrix I_r (steps 9 through 13).

Specifically, the algorithm decrees the following steps. (The steps as written include many parenthetical remarks—so many that some steps seem to consist more of parenthetical remarks than of actual algorithm. The remarks are unnecessary to execute the algorithm's steps as such. They are however necessary to explain and to justify the algorithm's steps to the reader.)

1. Begin by initializing

$$\begin{aligned}
 \tilde{P} &\leftarrow I, \\
 \tilde{D} &\leftarrow I, \\
 \tilde{L} &\leftarrow I, \\
 \tilde{U} &\leftarrow I, \\
 \tilde{I} &\leftarrow A, \\
 \tilde{K} &\leftarrow I, \\
 \tilde{S} &\leftarrow I, \\
 i &\leftarrow 1,
 \end{aligned}$$

where \tilde{I} holds the part of A remaining to be decomposed, where i is a row index, and where the others are the variable working matrices of (12.4). (The eventual goal will be to factor all of \tilde{I} away, leaving $\tilde{I} = I_r$, though the precise value of r will not be known until step 8. Since A by definition is a dimension-limited $m \times n$ matrix, one naturally need not store A beyond the $m \times n$ active region. What is less clear until one has read the whole algorithm, but nevertheless true, is that one also need not store the dimension-limited \tilde{I} beyond the $m \times n$ active region. The other six variable working matrices each have extended-operational form, but they also confine their activity to well defined regions: $m \times m$ for \tilde{P} , \tilde{D} , \tilde{L} and \tilde{U} ; $n \times n$ for \tilde{K} and \tilde{S} . One need store none of the matrices beyond these bounds.)

2. (Besides arriving at this point from step 1 above, the algorithm also reënters here from step 7 below. From step 1, $\tilde{I} = A$ and $\tilde{L} = I$, so this step 2 though logical seems unneeded. The need grows clear once one has read through step 7.) Observe that neither the i th row of \tilde{I} nor any row below it has an entry left of the i th column, that \tilde{I} is all-zero below-leftward of and directly leftward of (though not directly below) the *pivot* element \tilde{i}_{ii} .⁵ Observe also that above the i th row, the matrix has proper unit upper triangular form (§ 11.8). Regarding the other factors, notice that \tilde{L} enjoys the major partial unit triangular

⁵The notation \tilde{i}_{ii} looks interesting, but this is accidental. The \tilde{i} relates not to the doubled, subscribed index ii but to \tilde{I} . The notation \tilde{i}_{ii} thus means $[\tilde{I}]_{ii}$ —in other words, it means the current ii th element of the variable working matrix \tilde{I} .

form $L^{\{i-1\}}$ (§ 11.8.5) and that $\tilde{d}_{kk} = 1$ for all $k \geq i$. Pictorially,

$$\begin{aligned} \tilde{D} &= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & * & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & * & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & * & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \\ \tilde{L} = L^{\{i-1\}} &= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & * & 1 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & * & * & 1 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & * & * & * & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & * & * & * & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & * & * & * & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & * & * & * & 0 & 0 & 0 & 1 & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \\ \tilde{I} &= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & * & * & * & * & * & * & * & \cdots \\ \cdots & 0 & 1 & * & * & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 1 & * & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & * & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & * & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & * & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & * & * & * & * & * & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \end{aligned}$$

where the i th row and i th column are depicted at center.

3. Choose a nonzero element $\tilde{i}_{pq} \neq 0$ on or below the pivot row, where $p \geq i$ and $q \geq i$. (The easiest choice may simply be \tilde{i}_{ii} , where $p = q = i$, if $\tilde{i}_{ii} \neq 0$; but any nonzero element from the i th row downward can in general be chosen. Beginning students of the Gauss-Jordan or LU decomposition are conventionally taught to choose first the least possible q then the least possible p . When one has no reason to choose otherwise, that is as good a choice as any. There is however no actual need to choose so. In fact alternate choices can sometimes improve

practical numerical accuracy.^{6,7} Theoretically nonetheless, when doing exact arithmetic, the choice is quite arbitrary, so long as $\tilde{i}_{pq} \neq 0$.) If no nonzero element is available—if all remaining rows $p \geq i$ are now null—then skip directly to step 8.

4. Observing that (12.4) can be expanded to read

$$\begin{aligned} A &= \left(\tilde{P}T_{[p \leftrightarrow i]} \right) \left(T_{[p \leftrightarrow i]} \tilde{D}T_{[p \leftrightarrow i]} \right) \left(T_{[p \leftrightarrow i]} \tilde{L}T_{[p \leftrightarrow i]} \right) \left(T_{[p \leftrightarrow i]} \tilde{U}T_{[p \leftrightarrow i]} \right) \\ &\quad \times \left(T_{[p \leftrightarrow i]} \tilde{I}T_{[i \leftrightarrow q]} \right) \left(T_{[i \leftrightarrow q]} \tilde{K}T_{[i \leftrightarrow q]} \right) \left(T_{[i \leftrightarrow q]} \tilde{S} \right) \\ &= \left(\tilde{P}T_{[p \leftrightarrow i]} \right) \tilde{D} \left(T_{[p \leftrightarrow i]} \tilde{L}T_{[p \leftrightarrow i]} \right) \tilde{U} \\ &\quad \times \left(T_{[p \leftrightarrow i]} \tilde{I}T_{[i \leftrightarrow q]} \right) \tilde{K} \left(T_{[i \leftrightarrow q]} \tilde{S} \right), \end{aligned}$$

let

$$\begin{aligned} \tilde{P} &\leftarrow \tilde{P}T_{[p \leftrightarrow i]}, \\ \tilde{L} &\leftarrow T_{[p \leftrightarrow i]} \tilde{L}T_{[p \leftrightarrow i]}, \\ \tilde{I} &\leftarrow T_{[p \leftrightarrow i]} \tilde{I}T_{[i \leftrightarrow q]}, \\ \tilde{S} &\leftarrow T_{[i \leftrightarrow q]} \tilde{S}, \end{aligned}$$

⁶A typical Intel or AMD x86-class computer processor represents a C/C++ `double`-type floating-point number, $x = 2^p b$, in 0x40 bits of computer memory. Of the 0x40 bits, 0x34 are for the number's mantissa $2.0 \leq b < 4.0$ (not $1.0 \leq b < 2.0$ as one might expect), 0xB are for the number's exponent $-0x3FF \leq p \leq 0x3FE$, and one is for the number's \pm sign. (The mantissa's high-order bit, which is always 1, is implied not stored, thus is one neither of the 0x34 nor of the 0x40 bits.) The out-of-bounds exponents $p = -0x400$ and $p = 0x3FF$ serve specially respectively to encode 0 and ∞ . All this is standard computing practice. Such a floating-point representation is easily accurate enough for most practical purposes, but of course it is not generally exact. [22, § 1-4.2.2]

⁷The Gauss-Jordan's floating-point errors come mainly from dividing by small pivots. Such errors are naturally avoided by avoiding small pivots, at least until as late in the algorithm as possible. Smallness however is relative: a small pivot in a row and a column each populated by even smaller elements is unlikely to cause as much error as is a large pivot in a row and a column each populated by even larger elements.

To choose a pivot, any of several heuristics are reasonable. The following heuristic if programmed intelligently might not be too computationally expensive: Define the pivot-smallness metric

$$\tilde{\eta}_{pq}^2 \equiv \frac{2\tilde{i}_{pq}^* \tilde{i}_{pq}}{\sum_{p'=i}^m \tilde{i}_{p'q}^* \tilde{i}_{p'q} + \sum_{q'=i}^n \tilde{i}_{pq'}^* \tilde{i}_{pq'}}.$$

Choose the p and q of least $\tilde{\eta}_{pq}^2$. If two are equally least, then choose first the lesser column index q , then if necessary the lesser row index p .

(Though the step changes \tilde{L} , too, again it leaves \tilde{L} in the major partial unit triangular form $L^{\{i-1\}}$, because $i-1 < i$. Refer to Table 12.1.)

6. Observing that (12.4) can be expanded to read

$$\begin{aligned} A &= \tilde{P}\tilde{D} \left(\tilde{L}T_{\tilde{i}_{pi}[pi]} \right) \left(T_{-\tilde{i}_{pi}[pi]} \tilde{U}T_{\tilde{i}_{pi}[pi]} \right) \left(T_{-\tilde{i}_{pi}[pi]} \tilde{I} \right) \tilde{K}\tilde{S} \\ &= \tilde{P}\tilde{D} \left(\tilde{L}T_{\tilde{i}_{pi}[pi]} \right) \tilde{U} \left(T_{-\tilde{i}_{pi}[pi]} \tilde{I} \right) \tilde{K}\tilde{S}, \end{aligned}$$

clear \tilde{I} 's i th column below the pivot by letting

$$\begin{aligned} \tilde{L} &\leftarrow \left(\tilde{L} \right) \left(\prod_{p=i+1}^m T_{\tilde{i}_{pi}[pi]} \right), \\ \tilde{I} &\leftarrow \left(\prod_{p=i+1}^m T_{-\tilde{i}_{pi}[pi]} \right) \left(\tilde{I} \right). \end{aligned}$$

This forces $\tilde{i}_{ip} = 0$ for all $p > i$. It also fills in \tilde{L} 's i th column below the pivot, advancing that matrix from the $L^{\{i-1\}}$ form to the $L^{\{i\}}$ form. Pictorially,

$$\begin{aligned} \tilde{L} = L^{\{i\}} &= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & * & 1 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & * & * & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & * & * & * & 1 & 0 & 0 & 0 & \cdots \\ \cdots & * & * & * & * & 1 & 0 & 0 & \cdots \\ \cdots & * & * & * & * & 0 & 1 & 0 & \cdots \\ \cdots & * & * & * & * & 0 & 0 & 1 & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \\ \tilde{I} &= \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & * & * & * & * & * & * & \cdots \\ \cdots & 0 & 1 & * & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 1 & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 1 & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 0 & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 0 & * & * & * & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}. \end{aligned}$$

(Note that it is not necessary actually to apply the addition elementaries here one by one. Together they easily form an addition quasi-elementary $L_{[i]}$, thus can be applied all at once. See § 11.7.3.)

7. Increment

$$i \leftarrow i + 1.$$

Go back to step 2.

8. Decrement

$$i \leftarrow i - 1$$

to undo the last instance of step 7 (even if there never was an instance of step 7), thus letting i point to the matrix's last nonzero row. After decrementing, let the rank

$$r \equiv i.$$

Notice that, certainly, $r \leq m$ and $r \leq n$.

9. (Besides arriving at this point from step 8 above, the algorithm also reënters here from step 11 below.) If $i = 0$, then skip directly to step 12.

10. Observing that (12.4) can be expanded to read

$$A = \tilde{P}\tilde{D}\tilde{L} \left(\tilde{U}T_{i_{pi}[pi]} \right) \left(T_{-\tilde{i}_{pi}[pi]}\tilde{I} \right) \tilde{K}\tilde{S},$$

clear \tilde{I} 's i th column above the pivot by letting

$$\begin{aligned} \tilde{U} &\leftarrow \left(\tilde{U} \right) \left(\prod_{p=1}^{i-1} T_{\tilde{i}_{pi}[pi]} \right), \\ \tilde{I} &\leftarrow \left(\prod_{p=1}^{i-1} T_{-\tilde{i}_{pi}[pi]} \right) \left(\tilde{I} \right). \end{aligned}$$

This forces $\tilde{i}_{ip} = 0$ for all $p \neq i$. It also fills in \tilde{U} 's i th column above the pivot, advancing that matrix from the $U^{\{i+1\}}$ form to the $U^{\{i\}}$ form.

Pictorially,

$$\tilde{U} = U^{\{i\}} = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & 0 & 0 & * & * & * & * & \cdots \\ \cdots & 0 & 1 & 0 & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 1 & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 1 & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 1 & * & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 1 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

$$\tilde{I} = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & * & * & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 1 & * & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 1 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 1 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 1 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}.$$

(As in step 6, here again it is not necessary actually to apply the addition elementaries one by one. Together they easily form an addition quasialementary $U_{[i]}$. See § 11.7.3.)

11. Decrement $i \leftarrow i - 1$. Go back to step 9.
12. Notice that \tilde{I} now has the form of a rank- r identity matrix, except with $n - r$ extra columns dressing its right edge (often $r = n$ however; then there are no extra columns). Pictorially,

$$\tilde{I} = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & 1 & 0 & 0 & 0 & * & * & * & \cdots \\ \cdots & 0 & 1 & 0 & 0 & * & * & * & \cdots \\ \cdots & 0 & 0 & 1 & 0 & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 1 & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \cdots & 0 & 0 & 0 & 0 & 0 & 0 & 0 & \cdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}.$$

Observing that (12.4) can be expanded to read

$$A = \tilde{P}\tilde{D}\tilde{L}\tilde{U} \left(\tilde{I}T_{-\tilde{i}_{pq}[pq]} \right) \left(T_{\tilde{i}_{pq}[pq]}\tilde{K} \right) \tilde{S},$$

use the now conveniently elementarized columns of \tilde{I} 's main body to suppress the extra columns on its right edge by

$$\begin{aligned}\tilde{I} &\leftarrow (\tilde{I}) \left(\prod_{q=r+1}^n \prod_{p=1}^r T_{-\tilde{i}_{pq}[pq]} \right), \\ \tilde{K} &\leftarrow \left(\prod_{q=r+1}^n \prod_{p=1}^r T_{\tilde{i}_{pq}[pq]} \right) (\tilde{K}).\end{aligned}$$

(Actually, entering this step, it was that $\tilde{K} = I$, so in fact \tilde{K} becomes just the product above. As in steps 6 and 10, here again it is not necessary actually to apply the addition elementaries one by one. Together they easily form a parallel unit upper—not lower—triangular matrix $L_{\parallel}^{\{r\}T}$. See § 11.8.4.)

13. Notice now that $\tilde{I} = I_r$. Let

$$\begin{aligned}P &\equiv \tilde{P}, \\ D &\equiv \tilde{D}, \\ L &\equiv \tilde{L}, \\ U &\equiv \tilde{U}, \\ K &\equiv \tilde{K}, \\ S &\equiv \tilde{S}.\end{aligned}$$

End.

Never stalling, the algorithm cannot fail to achieve $\tilde{I} = I_r$ and thus a complete Gauss-Jordan decomposition of the form (12.2), though what value the rank r might turn out to have is not normally known to us in advance. (We have not yet proven, but shall in § 12.5, that the algorithm always produces the same I_r , the same rank $r \geq 0$, regardless of which pivots $\tilde{i}_{pq} \neq 0$ one happens to choose in step 3 along the way. We can safely ignore this unproven fact however for the immediate moment.)

12.3.4 Rank and independent rows

Observe that the Gauss-Jordan algorithm of § 12.3.3 operates always within the bounds of the original $m \times n$ matrix A . Therefore, necessarily,

$$\begin{aligned}r &\leq m, \\ r &\leq n.\end{aligned}\tag{12.5}$$

The rank r exceeds the number neither of the matrix's rows nor of its columns. This is unsurprising. Indeed the narrative of the algorithm's step 8 has already noticed the fact.

Observe also however that *the rank always fully reaches $r = m$ if the rows of the original matrix A are linearly independent*. The reason for this observation is that the rank can fall short, $r < m$, only if step 3 finds a null row $i \leq m$; but step 3 can find such a null row only if step 6 has created one (or if there were a null row in the original matrix A ; but according to § 12.1, such null rows never were linearly independent in the first place). How do we know that step 6 can never create a null row? We know this because the action of step 6 is to add multiples only of *current and earlier* pivot rows to rows in \tilde{I} which have not yet been on pivot.⁸ According to (12.1), such action has no power to cancel the independent rows it targets.

12.3.5 Inverting the factors

Inverting the six Gauss-Jordan factors is easy. Sections 11.7 and 11.8 have shown how. One need not however go even to that much trouble. Each of the six factors— P , D , L , U , K and S —is composed of a sequence $\coprod T$ of elementary operators. Each of the six inverse factors— P^{-1} , D^{-1} , L^{-1} , U^{-1} , K^{-1} and S^{-1} —is therefore composed of the *reverse* sequence $\coprod T^{-1}$ of *inverse* elementary operators. Refer to (11.41). If one merely records the sequence of elementaries used to build each of the six factors—if one reverses each sequence, inverts each elementary, and multiplies—then the six inverse factors result.

And, in fact, it isn't even that hard. One actually need not record the individual elementaries; one can invert, multiply and forget them in stream. This means starting the algorithm from step 1 with six extra variable work-

⁸If the truth of the sentence's assertion regarding the action of step 6 seems nonobvious, one can drown the assertion rigorously in symbols to prove it, but before going to that extreme consider: The action of steps 3 and 4 is to choose a pivot row $p \geq i$ and to shift it upward to the i th position. The action of step 6 then is to add multiples of the chosen pivot row downward only—that is, only to rows which have not yet been on pivot. This being so, steps 3 and 4 in the second iteration find no unmixed rows available to choose as second pivot, but find only rows which already include multiples of the first pivot row. Step 6 in the second iteration therefore adds multiples of the second pivot row downward, *which already include multiples of the first pivot row*. Step 6 in the i th iteration adds multiples of the i th pivot row downward, which already include multiples of the first through $(i-1)$ th. So it comes to pass that multiples only of current and earlier pivot rows are added to rows which have not yet been on pivot. To no row is ever added, directly or indirectly, a multiple of itself—until step 10, which does not belong to the algorithm's main loop and has nothing to do with the availability of nonzero rows to step 3.

ing matrices (besides the seven already there):

$$\begin{aligned}\tilde{P}^{-1} &\leftarrow I; \\ \tilde{D}^{-1} &\leftarrow I; \\ \tilde{L}^{-1} &\leftarrow I; \\ \tilde{U}^{-1} &\leftarrow I; \\ \tilde{K}^{-1} &\leftarrow I; \\ \tilde{S}^{-1} &\leftarrow I.\end{aligned}$$

(There is no \tilde{I}^{-1} , not because it would not be useful, but because its initial value would be⁹ $A^{-1(r)}$, unknown at algorithm's start.) Then, for each operation on any of \tilde{P} , \tilde{D} , \tilde{L} , \tilde{U} , \tilde{K} or \tilde{S} , one operates inversely on the corresponding inverse matrix. For example, in step 5,

$$\begin{aligned}\tilde{D} &\leftarrow \tilde{D}T_{\tilde{i}i}[i], & \tilde{D}^{-1} &\leftarrow T_{(1/\tilde{i}i)[i]}\tilde{D}^{-1}, \\ \tilde{L} &\leftarrow T_{(1/\tilde{i}i)[i]}\tilde{L}T_{\tilde{i}i}[i], & \tilde{L}^{-1} &\leftarrow T_{(1/\tilde{i}i)[i]}\tilde{L}^{-1}T_{\tilde{i}i}[i], \\ \tilde{I} &\leftarrow T_{(1/\tilde{i}i)[i]}\tilde{I}.\end{aligned}$$

With this simple extension, the algorithm yields all the factors not only of the Gauss-Jordan decomposition (12.2) but simultaneously also of the Gauss-Jordan's complementary form (12.3).

12.3.6 Truncating the factors

None of the six factors of (12.2) actually needs to retain its entire extended-operational form (§ 11.3.2). The four factors on the left, row operators, act wholly by their $m \times m$ squares; the two on the right, column operators, by their $n \times n$. Indeed, neither I_r nor A has anything but zeros outside the $m \times n$ rectangle, anyway, so there is nothing for the six operators to act upon beyond those bounds in any event. We can truncate all six operators to dimension-limited forms (§ 11.3.1) for this reason if we want.

To truncate the six operators formally, we left-multiply (12.2) by I_m and right-multiply it by I_n , obtaining

$$I_m A I_n = I_m P D L U I_r K S I_n.$$

According to § 11.3.6, the I_m and I_n respectively truncate rows and columns, actions which have no effect on A since it is already a dimension-limited $m \times n$

⁹Section 11.5 explains the notation.

matrix. By successive steps, then,

$$\begin{aligned} A &= I_m P D L U I_r K S I_n \\ &= I_m^7 P D L U I_r^2 K S I_n^3; \end{aligned}$$

and finally, by using (11.31) or (11.42) repeatedly,

$$A = (I_m P I_m)(I_m D I_m)(I_m L I_m)(I_m U I_r)(I_r K I_n)(I_n S I_n), \quad (12.6)$$

where the dimensionalities of the six factors on the equation's right side are respectively $m \times m$, $m \times m$, $m \times m$, $m \times r$, $r \times n$ and $n \times n$. Equation (12.6) expresses any dimension-limited rectangular matrix A as the product of six particularly simple, dimension-limited rectangular factors.

By similar reasoning from (12.2),

$$A = (I_m G_{>} I_r)(I_r G_{<} I_n), \quad (12.7)$$

where the dimensionalities of the two factors are $m \times r$ and $r \times n$.

The book will seldom point it out again explicitly, but one can straightforwardly truncate not only the Gauss-Jordan factors but most other factors and operators, too, by the method of this subsection.¹⁰

12.3.7 Properties of the factors

One would expect such neatly formed operators as the factors of the Gauss-Jordan to enjoy some useful special properties. Indeed they do. Table 12.2 lists a few. The table's properties formally come from (11.52) and Table 11.5;

¹⁰Comes the objection, "Black, why do you make it more complicated than it needs to be? For what reason must all your matrices have infinite dimensionality, anyway? They don't do it that way in my other linear algebra book."

It is a fair question. The answer is that this book is a book of applied mathematical theory; and theoretically in the author's view, infinite-dimensional matrices are significantly neater to handle. To append a null row or a null column to a dimension-limited matrix is to alter the matrix in no essential way, nor is there any real difference between $T_{5[21]}$ when it row-operates on a $3 \times p$ matrix and the same elementary when it row-operates on a $4 \times p$. The relevant theoretical constructs ought to reflect such insights. Hence infinite dimensionality.

Anyway, a matrix displaying an infinite field of zeros resembles a shipment delivering an infinite supply of nothing; one need not be too impressed with either. The two matrix forms of § 11.3 manifest the sense that a matrix can represent a linear transformation, whose rank matters; or a reversible row or column operation, whose rank does not. The extended-operational form, having infinite rank, serves the latter case. In either case, however, the dimensionality $m \times n$ of the matrix is a distraction. It is the rank r , if any, that counts.

Table 12.2: A few properties of the Gauss-Jordan factors.

$$\begin{aligned}
P^* &= P^{-1} = P^T \\
S^* &= S^{-1} = S^T \\
P^{-*} &= P = P^{-T} \\
S^{-*} &= S = S^{-T} \\
\frac{K + K^{-1}}{2} &= I \\
I_r K(I_n - I_r) &= K - I = I_r(K - I)(I_n - I_r) \\
I_r K^{-1}(I_n - I_r) &= K^{-1} - I = I_r(K^{-1} - I)(I_n - I_r) \\
(I - I_n)(K - I) &= 0 = (K - I)(I - I_n) \\
(I - I_n)(K^{-1} - I) &= 0 = (K^{-1} - I)(I - I_n)
\end{aligned}$$

but, if one firmly grasps the matrix forms involved and comprehends the notation (neither of which is trivial to do), if one understands that the operator $(I_n - I_r)$ is a truncator that selects columns $r + 1$ through n of the matrix it operates leftward upon, and if one sketches the relevant factors schematically with a pencil, then the properties are plainly seen without reference to Ch. 11 as such.

The table's properties regarding P and S express a general advantage all permutors share. The table's properties regarding K are admittedly less significant, included mostly only because § 13.3 will need them. Still, even the K properties are always true. They might find other uses.

Further properties of the several Gauss-Jordan factors can be gleaned from the respectively relevant subsections of §§ 11.7 and 11.8.

12.3.8 Marginalizing the factor I_n

If A happens to be a square, $n \times n$ matrix and if it develops that the rank $r = n$, then one can take advantage of (11.31) to rewrite the Gauss-Jordan decomposition (12.2) in the form

$$PDLUKSI_n = A = I_n PDLUKS, \quad (12.8)$$

thus marginalizing the factor I_n . This is to express the Gauss-Jordan solely in row operations or solely in column operations. It does not change the

algorithm and it does not alter the factors; it merely reorders the factors after the algorithm has determined them. It fails however if A is rectangular or $r < n$.

12.3.9 Decomposing an extended operator

Once one has derived the Gauss-Jordan decomposition, to extend it to decompose a reversible, $n \times n$ extended operator A (where per § 11.3.2 A outside the $n \times n$ active region resembles the infinite-dimensional identity matrix I) is trivial. One merely writes

$$A = PDLUKS,$$

wherein the I_r has become an I . Or, equivalently, one decomposes the $n \times n$ dimension-limited matrix $I_n A = I_n A I_n = A I_n$ as

$$A I_n = PDLU I_n K S = PDLUK S I_n,$$

from which, inasmuch as all the factors present but I_n are $n \times n$ extended operators, the preceding equation results.

One can decompose only reversible extended operators so. The Gauss-Jordan fails on irreversible extended operators, for which the rank of the truncated form $A I_n$ is $r < n$. See § 12.5.

This subsection's equations remain unnumbered because they say little new. Their only point, really, is that what an operator does outside some appropriately delimited active region is seldom interesting, because the vector on which the operator ultimately acts is probably null there in any event. In such a context it may not matter whether one truncates the operator. Indeed this was also the point of § 12.3.6 and, if you like, of (11.31), too.¹¹

¹¹If "it may not matter," as the narrative says, then one might just put all matrices in dimension-limited form. Many books do. To put them all in dimension-limited form however brings at least three effects the book you are reading prefers to avoid. First, it leaves shift-and-truncate operations hard to express cleanly (refer to §§ 11.3.6 and 11.9 and, as a typical example of the usage, eqn. 13.7). Second, it confuses the otherwise natural extension of discrete vectors into continuous functions. Third, it leaves one to consider the ranks of reversible operators like $T_{[1 \leftrightarrow 2]}$ that naturally should have no rank. The last of the three is arguably most significant: matrix rank is such an important attribute that one prefers to impute it only to those operators about which it actually says something interesting.

Nevertheless, the extended-operational matrix form is hardly more than a formality. All it says is that the extended operator unobtrusively leaves untouched anything it happens to find outside its operational domain, whereas a dimension-limited operator would have truncated whatever it found there. Since what is found outside the operational domain is often uninteresting, this may be a distinction without a difference, which one can safely ignore.

Regarding the present section as a whole, the Gauss-Jordan decomposition is a significant achievement. It is not the only matrix decomposition—further interesting decompositions include the Gram-Schmidt of § 13.10, the diagonal of § 14.6, the Schur of § 14.10 and the singular value of § 14.12, among others—but the Gauss-Jordan nonetheless reliably factors an arbitrary $m \times n$ matrix A , which we had not known how to handle very well, into a product of unit triangular matrices and quasiaelementaries, which we do. We shall put the Gauss-Jordan to good use in Ch. 13. However, before closing the present chapter we should like finally, squarely to define and to establish the concept of matrix rank, not only for I_r but for all matrices. To do that, we shall first need one more preliminary: the technique of vector replacement.

12.4 Vector replacement

Consider a set of $m + 1$ (not necessarily independent) vectors

$$\{\mathbf{u}, \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m\}.$$

As a definition, the *space* these vectors *address* consists of all linear combinations of the set's several vectors. That is, the space consists of all vectors \mathbf{b} formable as

$$\beta_o \mathbf{u} + \beta_1 \mathbf{a}_1 + \beta_2 \mathbf{a}_2 + \dots + \beta_m \mathbf{a}_m = \mathbf{b}. \quad (12.9)$$

Now consider a specific vector \mathbf{v} in the space,

$$\psi_o \mathbf{u} + \psi_1 \mathbf{a}_1 + \psi_2 \mathbf{a}_2 + \dots + \psi_m \mathbf{a}_m = \mathbf{v}, \quad (12.10)$$

for which

$$\psi_o \neq 0.$$

Solving (12.10) for \mathbf{u} , we find that

$$\frac{1}{\psi_o} \mathbf{v} - \frac{\psi_1}{\psi_o} \mathbf{a}_1 - \frac{\psi_2}{\psi_o} \mathbf{a}_2 - \dots - \frac{\psi_m}{\psi_o} \mathbf{a}_m = \mathbf{u}.$$

With the change of variables

$$\begin{aligned}\phi_o &\leftarrow \frac{1}{\psi_o}, \\ \phi_1 &\leftarrow -\frac{\psi_1}{\psi_o}, \\ \phi_2 &\leftarrow -\frac{\psi_2}{\psi_o}, \\ &\vdots \\ \phi_m &\leftarrow -\frac{\psi_m}{\psi_o},\end{aligned}$$

for which, quite symmetrically, it happens that

$$\begin{aligned}\psi_o &= \frac{1}{\phi_o}, \\ \psi_1 &= -\frac{\phi_1}{\phi_o}, \\ \psi_2 &= -\frac{\phi_2}{\phi_o}, \\ &\vdots \\ \psi_m &= -\frac{\phi_m}{\phi_o},\end{aligned}$$

the solution is

$$\phi_o \mathbf{v} + \phi_1 \mathbf{a}_1 + \phi_2 \mathbf{a}_2 + \cdots + \phi_m \mathbf{a}_m = \mathbf{u}. \quad (12.11)$$

Equation (12.11) has identical form to (12.10), only with the symbols $\mathbf{u} \leftrightarrow \mathbf{v}$ and $\psi \leftrightarrow \phi$ swapped. Since $\phi_o = 1/\psi_o$, assuming that ψ_o is finite it even appears that

$$\phi_o \neq 0;$$

so, the symmetry is complete. Table 12.3 summarizes.

Now further consider an arbitrary vector \mathbf{b} which lies in the space addressed by the vectors

$$\{\mathbf{u}, \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m\}.$$

Does the same \mathbf{b} also lie in the space addressed by the vectors

$$\{\mathbf{v}, \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m\}?$$

Table 12.3: The symmetrical equations of § 12.4.

$$\begin{array}{rcl}
\psi_o \mathbf{u} + \psi_1 \mathbf{a}_1 + \psi_2 \mathbf{a}_2 & & \phi_o \mathbf{v} + \phi_1 \mathbf{a}_1 + \phi_2 \mathbf{a}_2 \\
+ \cdots + \psi_m \mathbf{a}_m & = \mathbf{v} & + \cdots + \phi_m \mathbf{a}_m = \mathbf{u} \\
0 \neq \frac{1}{\psi_o} & = \phi_o & 0 \neq \frac{1}{\phi_o} = \psi_o \\
-\frac{\psi_1}{\psi_o} & = \phi_1 & -\frac{\phi_1}{\phi_o} = \psi_1 \\
-\frac{\psi_2}{\psi_o} & = \phi_2 & -\frac{\phi_2}{\phi_o} = \psi_2 \\
& \vdots & \vdots \\
-\frac{\psi_m}{\psi_o} & = \phi_m & -\frac{\phi_m}{\phi_o} = \psi_m
\end{array}$$

To show that it does, we substitute into (12.9) the expression for \mathbf{u} from (12.11), obtaining the form

$$(\beta_o)(\phi_o \mathbf{v} + \phi_1 \mathbf{a}_1 + \phi_2 \mathbf{a}_2 + \cdots + \phi_m \mathbf{a}_m) + \beta_1 \mathbf{a}_1 + \beta_2 \mathbf{a}_2 + \cdots + \beta_m \mathbf{a}_m = \mathbf{b}.$$

Collecting terms, this is

$$\beta_o \phi_o \mathbf{v} + (\beta_o \phi_1 + \beta_1) \mathbf{a}_1 + (\beta_o \phi_2 + \beta_2) \mathbf{a}_2 + \cdots + (\beta_o \phi_m + \beta_m) \mathbf{a}_m = \mathbf{b},$$

in which we see that, yes, \mathbf{b} does indeed also lie in the latter space. Naturally the problem's $\mathbf{u} \leftrightarrow \mathbf{v}$ symmetry then guarantees the converse, that an arbitrary vector \mathbf{b} which lies in the latter space also lies in the former. Therefore, a vector \mathbf{b} must lie in both spaces or neither, never just in one or the other. The two spaces are, in fact, one and the same.

This leads to the following useful conclusion. Given a set of vectors

$$\{\mathbf{u}, \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m\},$$

one can safely replace the \mathbf{u} by a new vector \mathbf{v} , obtaining the new set

$$\{\mathbf{v}, \mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m\},$$

provided that the replacement vector \mathbf{v} includes at least a little of the replaced vector \mathbf{u} ($\psi_o \neq 0$ in eqn. 12.10) and that \mathbf{v} is otherwise an honest

linear combination of the several vectors of the original set, untainted by foreign contribution. *Such vector replacement does not in any way alter the space addressed.* The new space is exactly the same as the old.

As a corollary, *if the vectors of the original set happen to be linearly independent (§ 12.1), then the vectors of the new set are linearly independent, too; for, if it were that*

$$\gamma_o \mathbf{v} + \gamma_1 \mathbf{a}_1 + \gamma_2 \mathbf{a}_2 + \cdots + \gamma_m \mathbf{a}_m = 0$$

for nontrivial γ_o and γ_k , then either $\gamma_o = 0$ —impossible since that would make the several \mathbf{a}_k themselves linearly dependent—or $\gamma_o \neq 0$, in which case \mathbf{v} would be a linear combination of the several \mathbf{a}_k alone. But if \mathbf{v} were a linear combination of the several \mathbf{a}_k alone, then (12.10) would still also explicitly make \mathbf{v} a linear combination of the same \mathbf{a}_k plus a nonzero multiple of \mathbf{u} . Yet both combinations cannot be, because according to § 12.1, two distinct combinations of independent vectors can never target the same \mathbf{v} . The contradiction proves false the assumption which gave rise to it: that the vectors of the new set were linearly dependent. Hence the vectors of the new set are equally as independent as the vectors of the old.

12.5 Rank

Sections 11.3.5 and 11.3.6 have introduced the rank- r identity matrix I_r , where the integer r is the number of ones the matrix has along its main diagonal. Other matrices have rank, too. Commonly, an $n \times n$ matrix has rank $r = n$, but consider the matrix

$$\begin{bmatrix} 5 & 1 & 6 \\ 3 & 6 & 9 \\ 2 & 4 & 6 \end{bmatrix}.$$

The third column of this matrix is the sum of the first and second columns. Also, the third row is just two-thirds the second. Either way, by columns or by rows, the matrix has only two independent vectors. The rank of this 3×3 matrix is not $r = 3$ but $r = 2$.

This section establishes properly the important concept of matrix *rank*. The section demonstrates that every matrix has a definite, unambiguous rank, and shows how this rank can be calculated.

To forestall potential confusion in the matter, we should immediately observe that—like the rest of this chapter but unlike some other parts of the book—this section explicitly trades in exact numbers. If a matrix element

here is 5, then it is exactly 5; if 0, then exactly 0. Many real-world matrices, of course—especially matrices populated by measured data—can never truly be exact, but that is not the point here. Here, the numbers are exact.¹²

12.5.1 A logical maneuver

In § 12.5.2 we shall execute a pretty logical maneuver, which one might name, “the end justifies the means.”¹³ When embedded within a larger logical construct as in § 12.5.2, the maneuver if unexpected can confuse, so this subsection is to prepare the reader to expect the maneuver.

The logical maneuver follows this basic pattern.

If P_1 then Q . If P_2 then Q . If P_3 then Q . Which of P_1 , P_2 and P_3 are true is not known, but what is known is that *at least one* of the three is true: P_1 or P_2 or P_3 . Therefore, although one can draw no valid conclusion regarding any one of the three predicates P_1 , P_2 or P_3 , one can still conclude that their common object Q is true.

One valid way to prove Q , then, would be to suppose P_1 and show that it led to Q , then alternately to suppose P_2 and show that it separately led to Q , then again to suppose P_3 and show that it also led to Q . The final step would be to show somehow that P_1 , P_2 and P_3 could not possibly all be false at once. Herein, the *means* is to assert several individually suspicious claims, none of which one actually means to prove. The *end* which justifies the means is the conclusion Q , which thereby one can and does prove.

It is a subtle maneuver. Once the reader feels that he grasps its logic, he can proceed to the next subsection where the maneuver is put to use.

¹²It is false to suppose that because applied mathematics *permits* imprecise quantities, like 3.0 ± 0.1 inches for the length of your thumb, it also requires them. On the contrary, the length of your thumb may indeed be 3.0 ± 0.1 inches, but surely no triangle has 3.0 ± 0.1 sides! A triangle has exactly three sides. The ratio of a circle’s circumference to its radius is exactly 2π . The author has exactly one brother. A construction contract might require the builder to finish within exactly 180 days (though the actual construction time might be an inexact $t = 172.6 \pm 0.2$ days), and so on. Exact quantities are every bit as valid in applied mathematics as imprecise ones are. Where the distinction matters, it is the applied mathematician’s responsibility to distinguish between the two kinds of quantity.

¹³The maneuver’s name rings a bit sinister, does it not? The author recommends no such maneuver in social or family life! Logically here however it helps the math.

12.5.2 The impossibility of identity-matrix promotion

Consider the matrix equation

$$AI_r B = I_s. \quad (12.12)$$

If $r \geq s$, then it is trivial to find matrices A and B for which (12.12) holds: $A = I_s = B$. If

$$r < s,$$

however, it is not so easy. In fact it is impossible. This subsection proves the impossibility. It shows that *one cannot by any row and column operations, reversible or otherwise, ever transform an identity matrix into another identity matrix of greater rank* (§ 11.3.5).

Equation (12.12) can be written in the form

$$(AI_r)B = I_s, \quad (12.13)$$

where, because I_r attacking from the right is the column truncation operator (§ 11.3.6), the product AI_r is a matrix with an unspecified number of rows but only r columns—or, more precisely, with no more than r nonzero columns. Viewed this way, per § 11.1.3, B operates on the r columns of AI_r to produce the s columns of I_s .

The r columns of AI_r are nothing more than the first through r th columns of A . Let the symbols $\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4, \mathbf{a}_5, \dots, \mathbf{a}_r$ denote these columns. The s columns of I_s , then, are nothing more than the elementary vectors $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3, \mathbf{e}_4, \mathbf{e}_5, \dots, \mathbf{e}_s$ (§ 11.3.7). The claim (12.13) makes is thus that the several vectors \mathbf{a}_k together address each of the several elementary vectors \mathbf{e}_j —that is, that a linear combination¹⁴

$$b_{1j}\mathbf{a}_1 + b_{2j}\mathbf{a}_2 + b_{3j}\mathbf{a}_3 + \dots + b_{rj}\mathbf{a}_r = \mathbf{e}_j \quad (12.14)$$

exists for each \mathbf{e}_j , $1 \leq j \leq s$.

The claim (12.14) will turn out to be false because there are too many \mathbf{e}_j , but to prove this, *we shall assume for the moment that the claim were true*. The proof then is by contradiction,¹⁵ and it runs as follows.

¹⁴Observe that unlike as in § 12.1, here we have not necessarily assumed that the several \mathbf{a}_k are linearly independent.

¹⁵As the reader will have observed by this point in the book, the technique—also called *reductio ad absurdum*—is the usual mathematical technique to prove impossibility. One assumes the falsehood to be true, then reasons toward a contradiction which proves the assumption false. Section 6.1.1 among others has already illustrated the technique, but the technique's use here is more sophisticated.

Consider the elementary vector \mathbf{e}_1 . For $j = 1$, (12.14) is

$$b_{11}\mathbf{a}_1 + b_{21}\mathbf{a}_2 + b_{31}\mathbf{a}_3 + \cdots + b_{r1}\mathbf{a}_r = \mathbf{e}_1,$$

which says that the elementary vector \mathbf{e}_1 is a linear combination of the several vectors

$$\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4, \mathbf{a}_5, \dots, \mathbf{a}_r\}.$$

Because \mathbf{e}_1 is a linear combination, according to § 12.4 one can safely replace any of the vectors in the set by \mathbf{e}_1 without altering the space addressed. For example, replacing \mathbf{a}_1 by \mathbf{e}_1 ,

$$\{\mathbf{e}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4, \mathbf{a}_5, \dots, \mathbf{a}_r\}.$$

The only restriction per § 12.4 is that \mathbf{e}_1 contain at least a little of the vector \mathbf{a}_k it replaces—that $b_{k1} \neq 0$. Of course there is no guarantee specifically that $b_{11} \neq 0$, so for \mathbf{e}_1 to replace \mathbf{a}_1 might not be allowed. However, inasmuch as \mathbf{e}_1 is nonzero, then according to (12.14) at least one of the several b_{k1} also is nonzero; and if b_{k1} is nonzero then \mathbf{e}_1 can replace \mathbf{a}_k . Some of the \mathbf{a}_k might indeed be forbidden, but never all; there is always at least one \mathbf{a}_k which \mathbf{e}_1 can replace. (For example, if \mathbf{a}_1 were forbidden because $b_{11} = 0$, then \mathbf{a}_3 might be available instead because $b_{31} \neq 0$. In this case the new set would be $\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{e}_1, \mathbf{a}_4, \mathbf{a}_5, \dots, \mathbf{a}_r\}$.)

Here is where the logical maneuver of § 12.5.1 comes in. The book to this point has established no general method to tell which of the several \mathbf{a}_k the elementary vector \mathbf{e}_1 actually contains (§ 13.2 gives the method, but that section depends logically on this one, so we cannot licitly appeal to it here). According to (12.14), the vector \mathbf{e}_1 might contain some of the several \mathbf{a}_k or all of them, but surely it contains at least one of them. Therefore, even though it is illegal to replace an \mathbf{a}_k by an \mathbf{e}_1 which contains none of it, even though we have no idea which of the several \mathbf{a}_k the vector \mathbf{e}_1 contains, even though replacing the wrong \mathbf{a}_k logically invalidates any conclusion which flows from the replacement, still we can proceed with the proof, provided that the false choice and the true choice lead ultimately alike to the same identical end. If they do, then all the logic requires is an assurance that there does exist at least one true choice, even if we remain ignorant as to which choice that is.

The claim (12.14) guarantees at least one true choice. Whether as the maneuver also demands, all the choices, true and false, lead ultimately alike to the same identical end remains to be determined.

Now consider the elementary vector \mathbf{e}_2 . According to (12.14), \mathbf{e}_2 lies in the space addressed by the original set

$$\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4, \mathbf{a}_5, \dots, \mathbf{a}_r\}.$$

Therefore as we have seen, \mathbf{e}_2 also lies in the space addressed by the new set

$$\{\mathbf{e}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4, \mathbf{a}_5, \dots, \mathbf{a}_r\}$$

(or $\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{e}_1, \mathbf{a}_4, \mathbf{a}_5, \dots, \mathbf{a}_r\}$, or whatever the new set happens to be). That is, not only do coefficients b_{k2} exist such that

$$b_{12}\mathbf{a}_1 + b_{22}\mathbf{a}_2 + b_{32}\mathbf{a}_3 + \dots + b_{r2}\mathbf{a}_r = \mathbf{e}_2,$$

but also coefficients β_{k2} exist such that

$$\beta_{12}\mathbf{e}_1 + \beta_{22}\mathbf{a}_2 + \beta_{32}\mathbf{a}_3 + \dots + \beta_{r2}\mathbf{a}_r = \mathbf{e}_2.$$

Again it is impossible for all the coefficients β_{k2} to be zero, but moreover, it is impossible for β_{12} to be the sole nonzero coefficient, for (as should seem plain to the reader who grasps the concept of the elementary vector, § 11.3.7) no elementary vector can ever be a linear combination of other elementary vectors alone! The linear combination which forms \mathbf{e}_2 evidently includes a nonzero multiple *of at least one of the remaining* \mathbf{a}_k . At least one of the β_{k2} attached to an \mathbf{a}_k (not β_{12} , which is attached to \mathbf{e}_1) must be nonzero. Therefore by the same reasoning as before, we now choose an \mathbf{a}_k with a nonzero coefficient $\beta_{k2} \neq 0$ and replace it by \mathbf{e}_2 , obtaining an even newer set of vectors like

$$\{\mathbf{e}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{e}_2, \mathbf{a}_5, \dots, \mathbf{a}_r\}.$$

This newer set addresses precisely the same space as the previous set, thus also as the original set.

And so it goes, replacing one \mathbf{a}_k by an \mathbf{e}_j at a time, until all the \mathbf{a}_k are gone and our set has become

$$\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3, \mathbf{e}_4, \mathbf{e}_5, \dots, \mathbf{e}_r\},$$

which, as we have reasoned, addresses precisely the same space as did the original set

$$\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4, \mathbf{a}_5, \dots, \mathbf{a}_r\}.$$

And this is the one identical end the maneuver of § 12.5.1 has demanded. All intermediate choices, true and false, ultimately lead to the single conclusion of this paragraph, which thereby is properly established.

The conclusion leaves us with a problem, however. There are more \mathbf{e}_j , $1 \leq j \leq s$, than there are \mathbf{a}_k , $1 \leq k \leq r$, because, as we have stipulated, $r < s$. Some elementary vectors \mathbf{e}_j , $r < j \leq s$, are evidently left over. Back at the beginning of the section, the claim (12.14) made was that

$$\{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_4, \mathbf{a}_5, \dots, \mathbf{a}_r\}$$

together addressed each of the several elementary vectors \mathbf{e}_j . But as we have seen, this amounts to a claim that

$$\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3, \mathbf{e}_4, \mathbf{e}_5, \dots, \mathbf{e}_r\}$$

together addressed each of the several elementary vectors \mathbf{e}_j . Plainly this is impossible with respect to the left-over \mathbf{e}_j , $r < j \leq s$. The contradiction proves false the claim which gave rise to it. The false claim: that the several \mathbf{a}_k , $1 \leq k \leq r$, addressed all the \mathbf{e}_j , $1 \leq j \leq s$, even when $r < s$.

Equation (12.13), which is just (12.12) written differently, asserts that B is a column operator which does precisely what we have just shown impossible: to combine the r columns of AI_r to yield the s columns of I_s , the latter of which are just the elementary vectors $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3, \mathbf{e}_4, \mathbf{e}_5, \dots, \mathbf{e}_s$. Hence finally we conclude that no matrices A and B exist which satisfy (12.12) when $r < s$. In other words, we conclude that *although row and column operations can demote identity matrices in rank, they can never promote them*. The promotion of identity matrices is impossible.

12.5.3 General matrix rank and its uniqueness

Step 8 of the Gauss-Jordan algorithm (§ 12.3.3) discovers a rank r for any matrix A . One should like to think that this rank r were a definite property of the matrix itself rather than some unreliable artifact of the algorithm, but until now we have lacked the background theory to prove it. Now we have the theory. Here is the proof.

The proof begins with a formal definition of the quantity whose uniqueness we are trying to prove.

- The rank r of an identity matrix I_r is the number of ones along its main diagonal. (This is from § 11.3.5.)
- The rank r of a general matrix A is the rank of an identity matrix I_r to which A can be reduced by *reversible* row and column operations.

A matrix A has rank r if and only if matrices $B_>$ and $B_<$ exist such that

$$\begin{aligned} B_>AB_< &= I_r, \\ A &= B_>^{-1}I_rB_<^{-1}, \\ B_>^{-1}B_> &= I = B_>B_>^{-1}, \\ B_<^{-1}B_< &= I = B_<B_<^{-1}. \end{aligned} \tag{12.15}$$

The question is whether in (12.15) only a single rank r is possible.

To answer the question, we suppose that another rank were possible, that A had not only rank r but also rank s . Then,

$$\begin{aligned} A &= B_>^{-1}I_rB_<^{-1}, \\ A &= G_>^{-1}I_sG_<^{-1}. \end{aligned}$$

Combining these equations,

$$B_>^{-1}I_rB_<^{-1} = G_>^{-1}I_sG_<^{-1}.$$

Solving first for I_r , then for I_s ,

$$\begin{aligned} (B_>G_>^{-1})I_s(G_<^{-1}B_<) &= I_r, \\ (G_>B_>^{-1})I_r(B_<^{-1}G_<) &= I_s. \end{aligned}$$

Were it that $r \neq s$, then one of these two equations would constitute the demotion of an identity matrix and the other, a promotion. But according to § 12.5.2 and its (12.12), promotion is impossible. Therefore $r \neq s$ is also impossible, and

$$r = s$$

is guaranteed. No matrix has two different ranks. *Matrix rank is unique.*

This finding has two immediate implications:

- Reversible row and/or column operations exist to change any matrix of rank r to *any other matrix* of the same rank. The reason is that, according to (12.15), reversible operations exist to change both matrices to I_r and back.
- No reversible operation can change a matrix's rank.

The discovery that every matrix has a single, unambiguous rank and the establishment of a failproof algorithm—the Gauss-Jordan—to ascertain that rank have not been easy to achieve, but they are important achievements

nonetheless, worth the effort thereto. The reason these achievements matter is that the mere dimensionality of a matrix is a chimerical measure of the matrix's true size—as for instance for the 3×3 example matrix at the head of the section. Matrix rank by contrast is an entirely solid, dependable measure. We shall rely on it often.

Section 12.5.8 comments further.

12.5.4 The full-rank matrix

According to (12.5), the rank r of a matrix can exceed the number neither of the matrix's rows nor of its columns. The greatest rank possible for an $m \times n$ matrix is the lesser of m and n . A *full-rank* matrix, then, is defined to be an $m \times n$ matrix with maximum rank $r = m$ or $r = n$ —or, if $m = n$, both. A matrix of less than full rank is a *degenerate* matrix.

Consider a tall $m \times n$ matrix C , $m \geq n$, one of whose n columns is a linear combination (§ 12.1) of the others. One could by definition target the dependent column with addition elementaries, using multiples of the other columns to wipe the dependent column out. Having zeroed the dependent column, one could then interchange it over to the matrix's extreme right, effectively throwing the column away, shrinking the matrix to $m \times (n - 1)$ dimensionality. Shrinking the matrix necessarily also shrinks the bound on the matrix's rank to $r \leq n - 1$ —which is to say, to $r < n$. But the shrink, done by reversible column operations, is itself reversible, by which § 12.5.3 binds the rank of the original, $m \times n$ matrix C likewise to $r < n$. The matrix C , one of whose columns is a linear combination of the others, is necessarily degenerate for this reason.

Now consider a tall matrix A with the same $m \times n$ dimensionality, but with a full n independent columns. The transpose A^T of such a matrix has a full n independent *rows*. One of the conclusions of § 12.3.4 was that a matrix of independent rows always has rank equal to the number of rows. Since A^T is such a matrix, its rank is a full $r = n$. But formally, what this says is that there exist operators $B_{<}^T$ and $B_{>}^T$ such that $I_n = B_{<}^T A^T B_{>}^T$, the transpose of which equation is $B_{>} A B_{<} = I_n$ —which in turn says that not only A^T , but also A itself, has full rank $r = n$.

Parallel reasoning rules the rows of broad matrices, $m \leq n$, of course. To square matrices, $m = n$, both lines of reasoning apply.

Gathering findings, we have that

- a tall $m \times n$ matrix, $m \geq n$, has full rank if and only if its columns are linearly independent;

- a broad $m \times n$ matrix, $m \leq n$, has full rank if and only if its rows are linearly independent;
- a square $n \times n$ matrix, $m = n$, has full rank if and only if its columns and/or its rows are linearly independent; and
- a square matrix has both independent columns and independent rows, or neither; never just one or the other.

To say that a matrix has *full column rank* is to say that it is tall or square and has full rank $r = n \leq m$. To say that a matrix has *full row rank* is to say that it is broad or square and has full rank $r = m \leq n$. Only a square matrix can have full column rank and full row rank at the same time, because a tall or broad matrix cannot but include, respectively, more columns or more rows than I_r .

12.5.5 Underdetermined and overdetermined linear systems (introduction)

The last paragraph of § 12.5.4 provokes yet further terminology. A linear system $A\mathbf{x} = \mathbf{b}$ is *underdetermined* if A lacks full column rank—that is, if $r < n$ —because inasmuch as some of A 's columns then depend linearly on the others such a system maps multiple n -element vectors \mathbf{x} to the same m -element vector \mathbf{b} , meaning that knowledge of \mathbf{b} does not suffice to determine \mathbf{x} uniquely. Complementarily, a linear system $A\mathbf{x} = \mathbf{b}$ is *overdetermined* if A lacks full row rank—that is, if $r < m$. If A lacks both, then the system is paradoxically both underdetermined and overdetermined and is thereby *degenerate*. If A happily has both, then the system is *exactly determined*.

Section 13.2 solves the exactly determined linear system. Section 13.4 solves the nonoverdetermined linear system. Section 13.6 analyzes the unsolvable overdetermined linear system among others. Further generalities await Ch. 13; but, regarding the overdetermined system specifically, the present subsection would observe at least the few following facts.

An overdetermined linear system $A\mathbf{x} = \mathbf{b}$ cannot have a solution for every possible m -element driving vector \mathbf{b} . The truth of this claim can be seen by decomposing the system's matrix A by Gauss-Jordan and then left-multiplying the decomposed system by $G_{>}^{-1}$ to reach the form

$$I_r G_{<} \mathbf{x} = G_{>}^{-1} \mathbf{b}.$$

If the m -element vector $\mathbf{c} \equiv G_{>}^{-1} \mathbf{b}$, then $I_r G_{<} \mathbf{x} = \mathbf{c}$, which is impossible unless the last $m - r$ elements of \mathbf{c} happen to be zero. But since $G_{>}$ is

invertible, each \mathbf{b} corresponds to a unique \mathbf{c} and vice versa; so, if \mathbf{b} is an unrestricted m -element vector then so also is \mathbf{c} , which verifies the claim.

Complementarily, *a nonoverdetermined linear system $A\mathbf{x} = \mathbf{b}$ does have a solution for every possible m -element driving vector \mathbf{b}* . This is so because in this case the last $m - r$ elements of \mathbf{c} do happen to be zero; or, better stated, because \mathbf{c} in this case has no nonzeros among its last $m - r$ elements, because it *has* no last $m - r$ elements, for the trivial reason that $r = m$.

It is an analytical error, and an easy one innocently to commit, to require that

$$A\mathbf{x} = \mathbf{b}$$

for unrestricted \mathbf{b} when A lacks full row rank. The error is easy to commit because the equation looks right, because such an equation is indeed valid over a broad domain of \mathbf{b} and might very well have been written correctly in that context, only not in the context of unrestricted \mathbf{b} . Analysis including such an error can lead to subtly absurd conclusions. It is never such an analytical error however to require that

$$A\mathbf{x} = 0$$

because, whatever other solutions such a system might have, it has at least the solution $\mathbf{x} = 0$.

12.5.6 The full-rank factorization

One sometimes finds dimension-limited matrices of less than full rank inconvenient to handle. However, every dimension-limited, $m \times n$ matrix of rank r can be expressed as the product of two full-rank matrices, one $m \times r$ and the other $r \times n$, both also of rank r :

$$A = BC. \tag{12.16}$$

The truncated Gauss-Jordan (12.7) constitutes one such *full-rank factorization*: $B = I_m G_{>} I_r$, $C = I_r G_{<} I_n$, good for any matrix. Other full-rank factorizations are possible, however, including among others the truncated Gram-Schmidt (13.54). The full-rank factorization is not unique.¹⁶

Of course, if an $m \times n$ matrix already has full rank $r = m$ or $r = n$, then the full-rank factorization is trivial: $A = I_m A$ or $A = A I_n$.

Section 13.6.4 uses the full-rank factorization.

¹⁶[2, § 3.3][33, “Moore-Penrose generalized inverse”]

12.5.7 Full column rank and the Gauss-Jordan factors K and S

The Gauss-Jordan decomposition (12.2),

$$A = PDLUI_rKS,$$

of a tall or square $m \times n$ matrix A of full column rank $r = n \leq m$ always finds the factor $K = I$, regardless of the pivots one chooses during the Gauss-Jordan algorithm's step 3. If one happens always to choose $q = i$ as pivot column then not only $K = I$ but $S = I$, too.

That $K = I$ is seen by the algorithm's step 12, which creates K . Step 12 nulls the spare columns $q > r$ that dress \tilde{I} 's right, but in this case \tilde{I} has only r columns and therefore has no spare columns to null. Hence step 12 does nothing and $K = I$.

That $S = I$ comes immediately of choosing $q = i$ for pivot column during each iterative instance of the algorithm's step 3. But, one must ask, can one choose so? What if column $q = i$ were unusable? That is, what if the only nonzero elements remaining in \tilde{I} 's i th column stood above the main diagonal, unavailable for step 4 to bring to pivot? Well, *were* it so, then one would indeed have to choose $q \neq i$ to swap the unusable column away rightward, but see: nothing in the algorithm later fills such a column's zeros with anything else—they remain zeros—so swapping the column away rightward could only delay the crisis. The column would remain unusable. Eventually the column would reappear on pivot when no usable column rightward remained available to swap it with, which contrary to our assumption would mean precisely that $r < n$. Such contradiction can only imply that if $r = n$ then no unusable column can ever appear. One need not swap. We conclude that though one might voluntarily choose $q \neq i$ during the algorithm's step 3, the algorithm cannot force one to do so if $r = n$. Yet if one always does choose $q = i$, as the full-column-rank matrix A evidently leaves one free to do, then indeed $S = I$.

Theoretically, the Gauss-Jordan decomposition (12.2) includes the factors K and S precisely to handle matrices with more columns than rank. Matrices of full column rank $r = n$, common in applications, by definition have no such problem. Therefore, the Gauss-Jordan decomposition theoretically needs no K or S for such matrices, which fact lets us abbreviate the decomposition for such matrices to read

$$A = PDLUI_n. \tag{12.17}$$

Observe however that just because one theoretically can set $S = I$ does not mean that one actually should. The column permutor S exists to be used, after all—especially numerically to avoid small pivots during early invocations of the algorithm’s step 5. Equation (12.17) is not mandatory but optional for a matrix A of full column rank (though still $r = n$ and thus $K = I$ for such a matrix, even when the unabbreviated eqn. 12.2 is used). There are however times when it is nice to know that one theoretically could, if doing exact arithmetic, set $S = I$ if one wanted to.

Since $PDLU$ acts as a row operator, (12.17) implies that each row of the full-rank matrix A lies in the space the rows of I_n address. This is obvious and boring, but interesting is the converse implication of (12.17)’s complementary form,

$$U^{-1}L^{-1}D^{-1}P^{-1}A = I_n,$$

that each row of I_n lies in the space the rows of A address. The rows of I_n and the rows of A evidently address the same space. One can moreover say the same of A ’s columns, since $B = A^T$ has full rank just as A does. In the whole, *if a matrix A is square and has full rank $r = n$, then A ’s columns together, A ’s rows together, I_n ’s columns together and I_n ’s rows together each address the same, complete n -dimensional space.*

12.5.8 The significance of rank uniqueness

The result of § 12.5.3, that matrix rank is unique, is an extremely important matrix theorem. It constitutes the chapter’s chief result, which we have spent so many pages to attain. Without this theorem, the very concept of matrix rank must remain in doubt, along with all that attends to the concept. The theorem is the rock upon which the general theory of the matrix is built.

The concept underlying the theorem promotes the useful sensibility that a matrix’s rank, much more than its mere dimensionality or the extent of its active region, represents the matrix’s true size. Dimensionality can deceive, after all. For example, the honest 2×2 matrix

$$\begin{bmatrix} 5 & 1 \\ 3 & 6 \end{bmatrix}$$

has two independent rows or, alternately, two independent columns, and, hence, rank $r = 2$. One can easily construct a phony 3×3 matrix from the honest 2×2 , however, simply by applying some 3×3 row and column

elementaries:

$$T_{(2/3)[32]} \begin{bmatrix} 5 & 1 \\ 3 & 6 \end{bmatrix} T_{1[13]} T_{1[23]} = \begin{bmatrix} 5 & 1 & 6 \\ 3 & 6 & 9 \\ 2 & 4 & 6 \end{bmatrix}.$$

The 3×3 matrix on the equation's right is the one we met at the head of the section. It *looks* like a rank-three matrix, but really has only two independent columns and two independent rows. Its true rank is $r = 2$. We have here caught a matrix impostor pretending to be bigger than it really is.¹⁷

Now, admittedly, adjectives like “honest” and “phony,” terms like “imposter,” are a bit hyperbolic. The last paragraph has used them to convey the subjective sense of the matter, but of course there is nothing mathematically improper or illegal about a matrix of less than full rank, so long as the true rank is correctly recognized. When one models a physical phenomenon by a set of equations, one sometimes is dismayed to discover that one of the equations, thought to be independent, is really just a useless combination of the others. This can happen in matrix work, too. The rank of a matrix helps one to recognize how many truly independent vectors, dimensions or equations one actually has available to work with, rather than how many seem available at first glance. That is the sense of matrix rank.

¹⁷An applied mathematician with some matrix experience actually probably recognizes this particular 3×3 matrix as a fraud on sight, but it is a very simple example. No one can just look at some arbitrary matrix and instantly perceive its true rank. Consider for instance the 5×5 matrix (in hexadecimal notation)

$$\begin{bmatrix} 12 & 9 & 3 & 1 & 0 \\ \frac{3}{2} & \frac{F}{2} & \frac{15}{2} & 2 & 12 \\ D & 9 & -19 & -\frac{E}{3} & -6 \\ -2 & 0 & 6 & 1 & 5 \\ 1 & -4 & 4 & 1 & -8 \end{bmatrix}.$$

As the reader can verify by the Gauss-Jordan algorithm, the matrix's rank is not $r = 5$ but $r = 4$.

Chapter 13

Inversion and orthonormalization

The undeniably tedious Chs. 11 and 12 have piled the matrix theory deep while affording scant practical reward. Building upon the two tedious chapters, this chapter brings the first rewarding matrix work.

One might be forgiven for forgetting after so many pages of abstract theory that the matrix afforded any reward or had any use at all. Uses however it has. Sections 11.1.1 and 12.5.5 have already broached¹ the matrix's most basic use, the primary subject of this chapter, to represent a system of m linear scalar equations in n unknowns neatly as

$$A\mathbf{x} = \mathbf{b}$$

and to solve the whole system at once by inverting the matrix A that characterizes it.

Now, before we go on, we want to confess that such a use alone, on the surface of it—though interesting—might not have justified the whole uncomfortable bulk of Chs. 11 and 12. We already knew how to solve a simultaneous system of linear scalar equations in principle without recourse to the formality of a matrix, after all, as in the last step to derive (3.9) as far back as Ch. 3. Why should we have suffered two bulky chapters, if only to prepare to do here something we already knew how to do?

The question is a fair one, but admits at least four answers. First, the matrix neatly solves a linear system not only for a particular driving vector \mathbf{b} but for all possible driving vectors \mathbf{b} at one stroke, as this chapter

¹The reader who has skipped Ch. 12 might at least review § 12.5.5.

explains. Second and yet more impressively, the matrix allows § 13.6 to introduce the *pseudoinverse* to approximate the solution to an unsolvable linear system and, moreover, to do so both optimally and efficiently, whereas such overdetermined systems arise commonly in applications. Third, to solve the linear system neatly is only the primary and most straightforward use of the matrix, not its only use: the even more interesting eigenvalue and its incidents await Ch. 14. Fourth, specific applications aside, one should never underestimate the blunt practical benefit of reducing an arbitrarily large grid of scalars to a single symbol A , which one can then manipulate by known algebraic rules. Most students first learning the matrix have wondered at this stage whether it were worth all the tedium; so, if the reader now wonders, then he stands in good company. The matrix finally begins to show its worth here.

The chapter opens in § 13.1 by inverting the square matrix to solve the exactly determined, $n \times n$ linear system in § 13.2. It continues in § 13.3 by computing the rectangular matrix's kernel to solve the nonoverdetermined, $m \times n$ linear system in § 13.4. In § 13.6, it brings forth the aforementioned pseudoinverse, which rightly approximates the solution to the unsolvable overdetermined linear system. After briefly revisiting the Newton-Raphson iteration in § 13.7, it concludes by introducing the concept and practice of vector orthonormalization in §§ 13.8 through 13.11.

13.1 Inverting the square matrix

Consider an $n \times n$ square matrix A of full rank $r = n$. Suppose that extended operators $G_>$, $G_<$, $G_>^{-1}$ and $G_<^{-1}$ can be found, each with an $n \times n$ active

region (§ 11.3.2), such that²

$$\begin{aligned} G_{>}^{-1}G_{>} &= I = G_{>}G_{>}^{-1}, \\ G_{<}^{-1}G_{<} &= I = G_{<}G_{<}^{-1}, \\ A &= G_{>}I_nG_{<}. \end{aligned} \tag{13.1}$$

Observing from (11.31) that

$$\begin{aligned} I_nA &= A = AI_n, \\ I_nG_{<}^{-1}G_{>}^{-1} &= G_{<}^{-1}I_nG_{>}^{-1} = G_{<}^{-1}G_{>}^{-1}I_n, \end{aligned}$$

we find by successive steps that

$$\begin{aligned} A &= G_{>}I_nG_{<}, \\ I_nA &= G_{>}G_{<}I_n, \\ G_{<}^{-1}G_{>}^{-1}I_nA &= I_n, \\ (G_{<}^{-1}I_nG_{>}^{-1})(A) &= I_n; \end{aligned}$$

or alternately that

$$\begin{aligned} A &= G_{>}I_nG_{<}, \\ AI_n &= I_nG_{>}G_{<}, \\ AI_nG_{<}^{-1}G_{>}^{-1} &= I_n, \\ (A)(G_{<}^{-1}I_nG_{>}^{-1}) &= I_n. \end{aligned}$$

²The symbology and associated terminology might disorient a reader who had skipped Chs. 11 and 12. In this book, the symbol I theoretically represents an $\infty \times \infty$ identity matrix, and $G_{>}$ and $G_{<}$ likewise run infinite strings of ones both directions out their main diagonals. When we say that “ P is an $n \times n$ permutor,” for example, we mean only that P is a permutor that resembles I outside its central, $n \times n$ active region—which is to say that P permutes only among the first n rows or columns. The symbol I_r contrarily represents an identity matrix of only r ones, though it too can be viewed as an $\infty \times \infty$ matrix with zeroes in the unused regions. If interpreted as an $\infty \times \infty$ matrix, the matrix A of the $m \times n$ system $A\mathbf{x} = \mathbf{b}$ has nonzero content only within the $m \times n$ rectangle.

None of this is complicated, really. Its purpose is merely to separate the essential features of a reversible operation like P from the dimensionality of the vector or matrix on which the operation happens to operate. The definitions do however necessarily, slightly diverge from definitions the reader may have been used to seeing in other books. In this book, one can legally multiply any two matrices, because all matrices are theoretically $\infty \times \infty$, anyway (though whether it makes any sense in a given circumstance to multiply mismatched matrices is another question; sometimes it does make sense, as in eqns. 13.24 and 14.49, but more often it does not—which naturally is why the other books tend to forbid such multiplication).

To the extent that the definitions confuse, the reader might briefly review the earlier chapters, especially § 11.3.

Either way, we have that

$$\begin{aligned} A^{-1}A &= I_n = AA^{-1}, \\ A^{-1} &\equiv G_{<}^{-1}I_nG_{>}^{-1}. \end{aligned} \tag{13.2}$$

Of course for this to work, $G_{>}$, $G_{<}$, $G_{>}^{-1}$ and $G_{<}^{-1}$ must exist, be known and honor $n \times n$ active regions, which might seem a practical hurdle. However, (12.2), (12.3) and the body of § 12.3 have shown exactly how to find just such a $G_{>}$, $G_{<}$, $G_{>}^{-1}$ and $G_{<}^{-1}$ for any square matrix A of full rank, without exception; so, there is no trouble here. The factors do exist, and indeed we know how to find them.

Equation (13.2) features the important matrix A^{-1} , the *rank- n inverse* of A . We have not yet much studied the rank- n inverse, but we have defined it in (11.49), where we gave it the fuller, nonstandard notation $A^{-1(n)}$. When naming the rank- n inverse in words one usually says simply, “the inverse,” because the rank is implied by the size of the square active region of the matrix inverted; but the rank- n inverse from (11.49) is not quite the infinite-dimensional inverse from (11.45), which is what $G_{>}^{-1}$ and $G_{<}^{-1}$ are. According to (13.2), the product of A^{-1} and A —or, written more fully, the product of $A^{-1(n)}$ and A —is, not I , but I_n .

Properties that emerge from (13.2) include the following.

- Like A , the rank- n inverse A^{-1} (more fully written $A^{-1(n)}$) too is an $n \times n$ square matrix of full rank $r = n$.
- Since A is square and has full rank (§ 12.5.4), its rows and, separately, its columns are linearly independent, so it has only the one, unique inverse A^{-1} . No other rank- n inverse of A exists.
- On the other hand, inasmuch as A is square and has full rank, it does per (13.2) indeed have an inverse A^{-1} . The rank- n inverse exists.
- If $B = A^{-1}$ then $B^{-1} = A$. That is, A is itself the rank- n inverse of A^{-1} . The matrices A and A^{-1} thus form an exclusive, reciprocal pair.
- If B is an $n \times n$ square matrix and either $BA = I_n$ or $AB = I_n$, then both equalities in fact hold; thus, $B = A^{-1}$. One can have neither equality without the other.
- Only a square, $n \times n$ matrix of full rank $r = n$ has a rank- n inverse. A matrix A' which is not square, or whose rank falls short of a full $r = n$, is not invertible in the rank- n sense of (13.2).

That A^{-1} is an $n \times n$ square matrix of full rank and that A is itself the inverse of A^{-1} proceed from the definition (13.2) of A^{-1} plus § 12.5.3's finding that reversible operations like $G_{>}^{-1}$ and $G_{<}^{-1}$ cannot change I_n 's rank. That the inverse exists is plain, inasmuch as the Gauss-Jordan decomposition plus (13.2) reliably calculate it. That the inverse is unique begins from § 12.5.4's observation that the columns (like the rows) of A are linearly independent because A is square and has full rank. From this beginning and the fact that $I_n = AA^{-1}$, it follows that $[A^{-1}]_{*1}$ represents³ the one and only possible combination of A 's columns which achieves \mathbf{e}_1 , that $[A^{-1}]_{*2}$ represents the one and only possible combination of A 's columns which achieves \mathbf{e}_2 , and so on through \mathbf{e}_n . One could observe likewise respecting the independent rows of A . Either way, A^{-1} is unique. Moreover, no other $n \times n$ matrix $B \neq A^{-1}$ satisfies *either* requirement of (13.2)—that $BA = I_n$ or that $AB = I_n$ —much less both.

It is not claimed that the matrix factors $G_{>}$ and $G_{<}$ themselves are unique, incidentally. On the contrary, many different pairs of matrix factors $G_{>}$ and $G_{<}$ can yield $A = G_{>}I_nG_{<}$, no less than that many different pairs of scalar factors $\gamma_{>}$ and $\gamma_{<}$ can yield $\alpha = \gamma_{>}1\gamma_{<}$. Though the Gauss-Jordan decomposition is a convenient means to $G_{>}$ and $G_{<}$, it is hardly the only means, and any proper $G_{>}$ and $G_{<}$ found by any means will serve so long as they satisfy (13.1). What are unique are not the factors but the A and A^{-1} they produce.

What of the degenerate $n \times n$ square matrix A' , of rank $r < n$? Rank promotion is impossible as §§ 12.5.2 and 12.5.3 have shown, so in the sense of (13.2) such a matrix has no inverse; for, if it had, then A'^{-1} would by definition represent a row or column operation which impossibly promoted A' to the full rank $r = n$ of I_n . Indeed, in that it has no inverse such a degenerate matrix closely resembles the scalar 0, which has no reciprocal. Mathematical convention owns a special name for a square matrix which is degenerate and thus has no inverse; it calls it a *singular* matrix.

And what of a rectangular matrix? Is it degenerate? Well, no, not exactly, not necessarily. The definitions of the present particular section are meant for square matrices; they do not neatly apply to nonsquare ones. Refer to §§ 12.5.3 and 12.5.4. However, appending the right number of null rows or columns to a nonsquare matrix does turn it into a degenerate square, in which case the preceding argument applies. See also §§ 12.5.5, 13.4 and 13.6.

³The notation $[A^{-1}]_{*j}$ means “the j th column of A^{-1} .” Refer to § 11.1.3.

13.2 The exactly determined linear system

Section 11.1.1 has shown how the single matrix equation

$$A\mathbf{x} = \mathbf{b} \tag{13.3}$$

concisely represents an entire simultaneous system of linear scalar equations. If the system has n scalar equations and n scalar unknowns, then the matrix A has square, $n \times n$ dimensionality. Furthermore, if the n scalar equations are independent of one another, then the rows of A are similarly independent, which gives A full rank and makes it invertible. Under these conditions, (13.3) and the corresponding system of linear scalar equations are solved simply by left-multiplying (13.3) by the A^{-1} of (13.2) and (13.1) to reach the famous formula

$$\mathbf{x} = A^{-1}\mathbf{b}. \tag{13.4}$$

Inverting the square matrix A of scalar coefficients, (13.4) concisely solves a simultaneous system of n linear scalar equations in n scalar unknowns. It is the classic motivational result of matrix theory.

If one omits first to prepare the theoretical ground sufficiently to support more advanced matrix work, then one can reach (13.4) with rather less effort than this book has done.⁴ As the chapter's introduction has observed, however, we shall soon meet additional interesting applications of the matrix which in any case require the theoretical ground to have been prepared. Equation (13.4) is only the first fruit of the effort.

⁴For motivational reasons, introductory, tutorial linear algebra textbooks like [20] and [27] rightly yet invariably invert the general square matrix of full rank much earlier, reaching (13.4) with less effort. The deferred price the student pays for the simpler-seeming approach of the tutorials is twofold. First, the student fails to develop the Gauss-Jordan decomposition properly, instead learning the less elegant but easier to grasp “row echelon form” of “Gaussian elimination” [20, Ch. 1][27, § 1.2]—which makes good matrix-arithmetic drill but leaves the student imperfectly prepared when the time comes to study kernels and eigensolutions or to read and write matrix-handling computer code. Second, in the long run the tutorials save no effort, because the student still must at some point develop the theory underlying matrix rank and supporting each of the several coincident properties of § 14.2. What the tutorials do is pedagogically necessary—it is how the writer first learned the matrix and probably how the reader first learned it, too—but it is appropriate to a tutorial, not to a study reference like this book.

In this book, where derivations prevail, the proper place to invert the general square matrix of full rank is here. Indeed, the inversion here goes smoothly, because Chs. 11 and 12 have laid under it a firm foundation upon which—and supplied it the right tools with which—to work.

Where the inverse does not exist, where the square matrix A is singular, the rows of the matrix are linearly dependent, meaning that the corresponding system actually contains fewer than n useful scalar equations. Depending on the value of the driving vector \mathbf{b} , the superfluous equations either merely reproduce or flatly contradict information the other equations already supply. Either way, no unique solution to a linear system described by a singular square matrix is possible—though a good approximate solution is given by the *pseudoinverse* of § 13.6. In the language of § 12.5.5, the singular square matrix characterizes a system that is both underdetermined and overdetermined, thus degenerate.

13.3 The kernel

If a matrix A has full column rank (§ 12.5.4), then the columns of A are linearly independent and

$$A\mathbf{x} = 0 \quad (13.5)$$

is impossible if $I_n\mathbf{x} \neq 0$. If the matrix however lacks full column rank then (13.5) is possible even if $I_n\mathbf{x} \neq 0$. In either case, any n -element \mathbf{x} (including $\mathbf{x} = 0$) that satisfies (13.5) belongs to the *kernel* of A .

Let A be an $m \times n$ matrix of rank r . A second matrix,⁵ A^K , minimally represents the kernel of A if and only if

- A^K has $n \times (n - r)$ dimensionality (which gives A^K tall rectangular form unless $r = 0$),
- A^K has full rank $n - r$ (that is, the columns of A^K are linearly independent, which gives A^K full column rank), and
- A^K satisfies the equation

$$AA^K = 0. \quad (13.6)$$

⁵The conventional mathematical notation for the kernel of A is $\ker\{A\}$, $\text{null}\{A\}$ or something nearly resembling one of the two—it seems to vary from editor to editor—which technically represent the kernel space itself, as opposed to the notation A^K which represents a matrix whose columns address the kernel space. This book deemphasizes the distinction and prefers the kernel matrix notation A^K .

If we were really precise, we might write not A^K but $A^{K(n)}$ to match the $A^{-1(r)}$ of (11.49). The abbreviated notation A^K is probably clear enough for most practical purposes, though, and surely more comprehensible to others who do not happen to have read this particular book.

The $n-r$ independent columns of the kernel matrix A^K address the complete space $\mathbf{x} = A^K \mathbf{a}$ of vectors in the kernel, where the $(n-r)$ -element vector \mathbf{a} can have any value. In symbols,

$$A\mathbf{x} = A(A^K \mathbf{a}) = (AA^K)\mathbf{a} = 0.$$

The definition does not pretend that the kernel matrix A^K is unique. Except when A has full column rank the kernel matrix is not unique; there are infinitely many kernel matrices A^K to choose from for a given matrix A . What is unique is not the kernel matrix but rather the space its columns address, and it is the latter space rather than A^K as such that is technically the kernel (if you forget and call A^K “a kernel,” though, you’ll be all right).

The *Gauss-Jordan kernel formula*⁶

$$A^K = S^{-1}K^{-1}H_r I_{n-r} = G_{<}^{-1}H_r I_{n-r} \quad (13.7)$$

gives a complete kernel A^K of A , where S^{-1} , K^{-1} and $G_{<}^{-1}$ are the factors their respective symbols indicate of the Gauss-Jordan decomposition’s complementary form (12.3) and H_r is the shift operator of § 11.9. Section 13.3.1 derives the formula, next.

13.3.1 The Gauss-Jordan kernel formula

To derive (13.7) is not easy. It begins from the statement of the linear system

$$A\mathbf{x} = \mathbf{b}, \text{ where } \mathbf{b} = 0 \text{ or } r = m, \text{ or both;} \quad (13.8)$$

and where \mathbf{b} and \mathbf{x} are respectively m - and n -element vectors and A is an $m \times n$ matrix of rank r . This statement is broader than (13.5) requires but it serves § 13.4, too; so, for the moment, for generality’s sake, we leave \mathbf{b} unspecified but by the given proviso. Gauss-Jordan factoring A , by successive steps,

$$\begin{aligned} G_{>} I_r K S \mathbf{x} &= \mathbf{b}, \\ I_r K S \mathbf{x} &= G_{>}^{-1} \mathbf{b}, \\ I_r (K - I) S \mathbf{x} + I_r S \mathbf{x} &= G_{>}^{-1} \mathbf{b}. \end{aligned}$$

Applying an identity from Table 12.2 on page 300,

$$I_r K (I_n - I_r) S \mathbf{x} + I_r S \mathbf{x} = G_{>}^{-1} \mathbf{b}.$$

⁶The name *Gauss-Jordan kernel formula* is not standard as far as the writer is aware, but we would like a name for (13.7). This name seems as fitting as any.

Rearranging terms,

$$I_r S\mathbf{x} = G_{>}^{-1} \mathbf{b} - I_r K(I_n - I_r) S\mathbf{x}. \quad (13.9)$$

Equation (13.9) is interesting. It has $S\mathbf{x}$ on both sides, where $S\mathbf{x}$ is the vector \mathbf{x} with elements reordered in some particular way. The equation has however on the left only $I_r S\mathbf{x}$, which is the first r elements of $S\mathbf{x}$; and on the right only $(I_n - I_r) S\mathbf{x}$, which is the remaining $n - r$ elements.⁷ No element of $S\mathbf{x}$ appears on both sides. Naturally this is no accident; we have (probably after some trial and error not recorded here) planned the steps leading to (13.9) to achieve precisely this effect. Equation (13.9) implies *that one can choose the last $n - r$ elements of $S\mathbf{x}$ freely, but that the choice then determines the first r elements.*

The implication is significant. To express the implication more clearly we can rewrite (13.9) in the improved form

$$\begin{aligned} \mathbf{f} &= G_{>}^{-1} \mathbf{b} - I_r K H_r \mathbf{a}, \\ S\mathbf{x} &= \begin{bmatrix} \mathbf{f} \\ \mathbf{a} \end{bmatrix} = \mathbf{f} + H_r \mathbf{a}, \\ \mathbf{f} &\equiv I_r S\mathbf{x}, \\ \mathbf{a} &\equiv H_{-r}(I_n - I_r) S\mathbf{x}, \end{aligned} \quad (13.10)$$

where \mathbf{a} represents the $n - r$ free elements of $S\mathbf{x}$ and \mathbf{f} represents the r dependent elements. This makes \mathbf{f} and thereby also \mathbf{x} functions of the free parameter \mathbf{a} and the driving vector \mathbf{b} :

$$\begin{aligned} \mathbf{f}(\mathbf{a}, \mathbf{b}) &= G_{>}^{-1} \mathbf{b} - I_r K H_r \mathbf{a}, \\ S\mathbf{x}(\mathbf{a}, \mathbf{b}) &= \begin{bmatrix} \mathbf{f}(\mathbf{a}, \mathbf{b}) \\ \mathbf{a} \end{bmatrix} = \mathbf{f}(\mathbf{a}, \mathbf{b}) + H_r \mathbf{a}. \end{aligned} \quad (13.11)$$

If $\mathbf{b} = 0$ as (13.5) requires, then

$$\begin{aligned} \mathbf{f}(\mathbf{a}, 0) &= -I_r K H_r \mathbf{a}, \\ S\mathbf{x}(\mathbf{a}, 0) &= \begin{bmatrix} \mathbf{f}(\mathbf{a}, 0) \\ \mathbf{a} \end{bmatrix} = \mathbf{f}(\mathbf{a}, 0) + H_r \mathbf{a}. \end{aligned}$$

Substituting the first line into the second,

$$S\mathbf{x}(\mathbf{a}, 0) = (I - I_r K) H_r \mathbf{a}. \quad (13.12)$$

⁷Notice how we now associate the factor $(I_n - I_r)$ rightward as a row truncator, though it had first entered acting leftward as a column truncator. The flexibility to reassociate operators in such a way is one of many good reasons Chs. 11 and 12 have gone to such considerable trouble to develop the basic theory of the matrix.

In the event that $\mathbf{a} = \mathbf{e}_j$, where $1 \leq j \leq n - r$,

$$S\mathbf{x}(\mathbf{e}_j, 0) = (I - I_r K)H_r \mathbf{e}_j.$$

For all the \mathbf{e}_j at once,

$$S\mathbf{x}(I_{n-r}, 0) = (I - I_r K)H_r I_{n-r}.$$

But if all the \mathbf{e}_j at once, the columns of I_{n-r} , exactly address the domain of \mathbf{a} , then the columns of $\mathbf{x}(I_{n-r}, 0)$ likewise exactly address the range of $\mathbf{x}(\mathbf{a}, 0)$. Equation (13.6) has already named this range A^K , by which⁸

$$SA^K = (I - I_r K)H_r I_{n-r}. \quad (13.13)$$

Left-multiplying by $S^{-1} = S^* = S^T$ produces the alternate kernel formula

$$A^K = S^{-1}(I - I_r K)H_r I_{n-r}. \quad (13.14)$$

The alternate kernel formula (13.14) is correct but not as simple as it could be. By the identity (11.76), (13.13) is

$$\begin{aligned} SA^K &= (I - I_r K)(I_n - I_r)H_r \\ &= [(I_n - I_r) - I_r K(I_n - I_r)]H_r \\ &= [(I_n - I_r) - (K - I)]H_r, \end{aligned} \quad (13.15)$$

⁸These are difficult steps. How does one justify replacing \mathbf{a} by \mathbf{e}_j , then \mathbf{e}_j by I_{n-r} , then \mathbf{x} by A^K ? One justifies them in that the columns of I_{n-r} are the several \mathbf{e}_j , of which any $(n - r)$ -element vector \mathbf{a} can be constructed as the linear combination

$$\mathbf{a} = I_{n-r}\mathbf{a} = [\mathbf{e}_1 \quad \mathbf{e}_2 \quad \mathbf{e}_3 \quad \cdots \quad \mathbf{e}_{n-r}] \mathbf{a} = \sum_{j=1}^{n-r} a_j \mathbf{e}_j$$

weighted by the elements of \mathbf{a} . Seen from one perspective, this seems trivial; from another perspective, baffling; until one grasps what is really going on here.

The idea is that if we can solve the problem for each elementary vector \mathbf{e}_j —that is, in aggregate, if we can solve the problem for the identity matrix I_{n-r} —then we shall implicitly have solved it for every \mathbf{a} because \mathbf{a} is a weighted combination of the \mathbf{e}_j and the whole problem is linear. The solution

$$\mathbf{x} = A^K \mathbf{a}$$

for a given choice of \mathbf{a} becomes a weighted combination of the solutions for each \mathbf{e}_j , with the elements of \mathbf{a} again as the weights. And what are the solutions for each \mathbf{e}_j ? Answer: the corresponding columns of A^K , which by definition are the independent values of \mathbf{x} that cause $\mathbf{b} = 0$.

where we have used Table 12.2 again in the last step. How to proceed symbolically from (13.15) is not obvious, but if one sketches the matrices of (13.15) schematically with a pencil, and if one remembers that K^{-1} is just K with elements off the main diagonal negated, then it appears that

$$SA^K = K^{-1}H_r I_{n-r}. \quad (13.16)$$

The appearance is not entirely convincing,⁹ but (13.16) though unproven still helps because it posits a hypothesis toward which to target the analysis.

Two variations on the identities of Table 12.2 also help. First, from the identity that

$$\frac{K + K^{-1}}{2} = I,$$

we have that

$$K - I = I - K^{-1}. \quad (13.17)$$

Second, right-multiplying by I_r the identity that

$$I_r K^{-1}(I_n - I_r) = K^{-1} - I$$

and canceling terms, we have that

$$K^{-1}I_r = I_r \quad (13.18)$$

(which actually is pretty obvious if you think about it, since all of K 's interesting content lies by construction right of its r th column). Now we have enough to go on with. Substituting (13.17) and (13.18) into (13.15) yields

$$SA^K = [(I_n - K^{-1}I_r) - (I - K^{-1})]H_r.$$

Adding $0 = K^{-1}I_n H_r - K^{-1}I_n H_r$ and rearranging terms,

$$SA^K = K^{-1}(I_n - I_r)H_r + [K^{-1} - K^{-1}I_n - I + I_n]H_r.$$

Factoring,

$$SA^K = K^{-1}(I_n - I_r)H_r + [(K^{-1} - I)(I - I_n)]H_r.$$

According to Table 12.2, the quantity in square brackets is zero, so

$$SA^K = K^{-1}(I_n - I_r)H_r,$$

⁹Well, no, actually, the appearance pretty much is entirely convincing, but let us finish the proof symbolically nonetheless.

which, considering that the identity (11.76) has that $(I_n - I_r)H_r = H_r I_{n-r}$, proves (13.16). The final step is to left-multiply (13.16) by $S^{-1} = S^* = S^T$, reaching (13.7) that was to be derived.

One would like to feel sure that the columns of (13.7)'s A^K actually addressed the whole kernel space of A rather than only part. One would further like to feel sure that A^K had no redundant columns; that is, that it had full rank. Moreover, the definition of A^K in the section's introduction demands both of these features. In general such features would be hard to establish, but here the factors conveniently are Gauss-Jordan factors. Regarding the whole kernel space, A^K addresses it because A^K comes from all \mathbf{a} . Regarding redundancy, A^K lacks it because SA^K lacks it, and SA^K lacks it because according to (13.13) the last rows of SA^K are $H_r I_{n-r}$. So, in fact, (13.7) has both features and does fit the definition.

13.3.2 Converting between kernel matrices

If C is a reversible $(n-r) \times (n-r)$ operator by which we right-multiply (13.6), then the matrix

$$A'^K = A^K C \quad (13.19)$$

like A^K evidently represents the kernel of A :

$$AA'^K = A(A^K C) = (AA^K)C = 0.$$

Indeed this makes sense: because the columns of $A^K C$ address the same space the columns of A^K address, the two matrices necessarily represent the same underlying kernel. Moreover, *some* C exists to convert A^K into every alternate kernel matrix A'^K of A . We know this because § 12.4 lets one replace the columns of A^K with those of A'^K , reversibly, one column at a time, without altering the space addressed. (It might not let one replace the columns in sequence, but if out of sequence then a reversible permutation at the end corrects the order. Refer to §§ 12.5.1 and 12.5.2 for the pattern by which this is done.)

The orthonormalizing column operator R^{-1} of (13.52) below incidentally tends to make a good choice for C .

13.3.3 The degree of freedom

A slightly vague but extraordinarily useful concept has emerged in this section, worth pausing briefly to appreciate. The concept is the concept of the *degree of freedom*.

In short, a degree of freedom is a parameter one remains free to determine within some continuous domain. Napoleon's artillerist¹⁰ might for example have enjoyed as many as six degrees of freedom in firing a cannonball: two in where he chose to set up his cannon (one degree in north-south position, one in east-west); two in aim (azimuth and elevation); one in muzzle velocity (governed by the quantity of gunpowder used to propel the ball); and one in time. A seventh potential degree of freedom, the height from which the artillerist fires, is of course restricted by the lay of the land: the artillerist can fire from a high place only if the place he has chosen to fire from happens to be up on a hill (Napoleon had no flying cannon). Yet even among the six remaining degrees of freedom, the artillerist might find some impractical to exercise. The artillerist probably preloads the cannon always with a standard charge of gunpowder because, when he finds his target in the field, he cannot spare the time to unload the cannon and alter the charge: this costs one degree of freedom. Likewise, the artillerist must limber up the cannon and hitch it to a horse to shift it to better ground; for this too he cannot spare time in the heat of battle: this costs two degrees. And Napoleon might yell, "Fire!" canceling the time degree as well. Two degrees of freedom remain to the artillerist; but, since exactly two degrees are needed to hit some particular target on the battlefield, the two are enough.

Now consider what happens if the artillerist loses one of his last two remaining degrees of freedom. Maybe the cannon's carriage wheel is broken and the artillerist can no longer turn the cannon; that is, he can still choose firing elevation but no longer azimuth. In such a strait to hit some particular target on the battlefield, the artillerist needs somehow to recover another degree of freedom, for he needs two but has only one. If he disregards Napoleon's order, "Fire!" (maybe not a wise thing to do, but, anyway, ...) and waits for the target to traverse the cannon's fixed line of fire, then he can still hope to hit even with the broken carriage wheel; for could he choose neither azimuth nor the moment to fire, then he would almost surely miss.

Some apparent degrees of freedom are not real. For example, muzzle velocity gives the artillerist little control firing elevation does not also give. Other degrees of freedom are nonlinear in effect: a certain firing elevation gives maximum range; nearer targets can be hit by firing either higher or lower at the artillerist's discretion. On the other hand, too much gunpowder might break the cannon.

¹⁰The author, who has never fired an artillery piece (unless an arrow from a Boy Scout bow counts), invites any real artillerist among the readership to write in to improve the example.

All of this is hard to generalize in unambiguous mathematical terms, but the count of the degrees of freedom in a system is of high conceptual importance to the engineer nonetheless. Basically, the count captures the idea that in order to control n output variables of some system, one needs at least n independent input variables. The n may possibly for various reasons still not suffice—it might be wise in some cases to allow $n + 1$ or $n + 2$ —but in no event will fewer than n do. Engineers of all kinds think in this way: an aeronautical engineer knows in advance that an airplane needs at least n ailerons, rudders and other control surfaces for the pilot adequately to control the airplane; an electrical engineer knows in advance that a circuit needs at least n potentiometers for the technician adequately to tune the circuit; and so on.

In geometry, a line brings a single degree of freedom. A plane brings two. A point brings none. If the line bends and turns like a mountain road, it still brings a single degree of freedom. And if the road reaches an intersection? Answer: still one degree. A degree of freedom has some continuous nature, not merely a discrete choice to turn left or right. On the other hand, a swimmer in a swimming pool enjoys three degrees of freedom (up-down, north-south, east-west) even though his domain in any of the three is limited to the small volume of the pool. The driver on the mountain road cannot claim a second degree of freedom at the mountain intersection (he can indeed claim a choice, but the choice being discrete lacks the proper character of a degree of freedom), but he might plausibly claim a second degree of freedom upon reaching the city, where the web or grid of streets is dense enough to approximate access to any point on the city's surface. Just how many streets it takes to turn the driver's "line" experience into a "plane" experience is a matter for the mathematician's discretion.

Reviewing (13.11), we find $n - r$ degrees of freedom in the general underdetermined linear system, represented by the $n - r$ free elements of \mathbf{a} . If the underdetermined system is not also overdetermined, if it is nondegenerate such that $r = m$, then it is guaranteed to have a family of solutions \mathbf{x} . This family is the topic of the next section.

13.4 The nonoverdetermined linear system

The exactly determined linear system of § 13.2 is common, but also common is the more general, nonoverdetermined linear system

$$A\mathbf{x} = \mathbf{b}, \tag{13.20}$$

in which \mathbf{b} is a known, m -element vector; \mathbf{x} is an unknown, n -element vector; and A is a square or broad, $m \times n$ matrix of full row rank (§ 12.5.4)

$$r = m \leq n. \quad (13.21)$$

Except in the exactly determined edge case $r = m = n$ of § 13.2, the nonoverdetermined linear system has no unique solution but rather a family of solutions. This section delineates the family.

13.4.1 Particular and homogeneous solutions

The nonoverdetermined linear system (13.20) by definition admits more than one solution \mathbf{x} for a given driving vector \mathbf{b} . Such a system is hard to solve all at once, though, so we prefer to split the system as

$$\begin{aligned} A\mathbf{x}_1 &= \mathbf{b}, \\ A(A^K\mathbf{a}) &= 0, \\ \mathbf{x} &= \mathbf{x}_1 + A^K\mathbf{a}, \end{aligned} \quad (13.22)$$

which, when the second line is added to the first and the third is substituted, makes the whole form (13.20). Splitting the system does not change it, but it does let us treat the system's first and second lines in (13.22) separately. In the split form, the symbol \mathbf{x}_1 represents any one n -element vector that happens to satisfy the form's first line—many are possible; the mathematician just picks one—and is called *a particular solution* of (13.20). The $(n - r)$ -element vector \mathbf{a} remains unspecified, whereupon $A^K\mathbf{a}$ represents the complete family of n -element vectors that satisfy the form's second line. The family of vectors expressible as $A^K\mathbf{a}$ is called *the homogeneous solution* of (13.20).

Notice the italicized articles *a* and *the*.

The Gauss-Jordan kernel formula (13.7) has given us A^K and thereby the homogeneous solution, which renders the analysis of (13.20) already half done. To complete the analysis, it remains in § 13.4.2 to find a particular solution.

13.4.2 A particular solution

Any particular solution will do. Equation (13.11) has that

$$\begin{aligned} \mathbf{f}(\mathbf{a}, \mathbf{b}) &= G_{>}^{-1}\mathbf{b} - I_r K H_r \mathbf{a}, \\ (S) [\mathbf{x}_1(\mathbf{a}, \mathbf{b}) + A^K\mathbf{a}] &= \begin{bmatrix} \mathbf{f}(\mathbf{a}, \mathbf{b}) \\ \mathbf{a} \end{bmatrix} = \mathbf{f}(\mathbf{a}, \mathbf{b}) + H_r \mathbf{a}, \end{aligned}$$

where we have substituted the last line of (13.22) for \mathbf{x} . This holds for any \mathbf{a} and \mathbf{b} . We are not free to choose the driving vector \mathbf{b} , but since we need only one particular solution, \mathbf{a} can be anything we want. Why not

$$\mathbf{a} = \mathbf{0}?$$

Then

$$\begin{aligned} \mathbf{f}(0, \mathbf{b}) &= G_{>}^{-1} \mathbf{b}, \\ S\mathbf{x}_1(0, \mathbf{b}) &= \begin{bmatrix} \mathbf{f}(0, \mathbf{b}) \\ 0 \end{bmatrix} = \mathbf{f}(0, \mathbf{b}). \end{aligned}$$

That is,

$$\mathbf{x}_1 = S^{-1} G_{>}^{-1} \mathbf{b}. \quad (13.23)$$

13.4.3 The general solution

Assembling (13.7), (13.22) and (13.23) yields the general solution

$$\mathbf{x} = S^{-1}(G_{>}^{-1} \mathbf{b} + K^{-1} H_r I_{n-r} \mathbf{a}) \quad (13.24)$$

to the nonoverdetermined linear system (13.20).

In exact arithmetic (13.24) solves the nonoverdetermined linear system in theory exactly. Of course, practical calculations are usually done in limited precision, in which compounded rounding error in the last bit eventually disrupts (13.24) for matrices larger than some moderately large size. Avoiding unduly small pivots early in the Gauss-Jordan extends (13.24)'s reach to larger matrices, and for yet larger matrices a bewildering variety of more sophisticated techniques exists to mitigate the problem, which can be vexing because the problem arises even when the matrix A is exactly known. Equation (13.24) is useful and correct, but one should at least be aware that it can in practice lose floating-point accuracy when the matrix it attacks grows too large. (It can also lose accuracy when the matrix's rows are almost dependent, but that is more the fault of the matrix than of the formula. See § 14.8, which addresses a related problem.)

13.5 The residual

Equations (13.2) and (13.4) solve the exactly determined linear system $A\mathbf{x} = \mathbf{b}$. Equation (13.24) broadens the solution to include the nonoverdetermined linear system. None of those equations however can handle the

overdetermined linear system, because for general \mathbf{b} the overdetermined linear system

$$A\mathbf{x} \approx \mathbf{b} \quad (13.25)$$

has no exact solution. (See § 12.5.5 for the definitions of *underdetermined*, *overdetermined*, etc.)

One is tempted to declare the overdetermined system uninteresting because it has no solution and to leave the matter there, but this would be a serious mistake. In fact the overdetermined system is especially interesting, and the more so because it arises so frequently in applications. One seldom trusts a minimal set of data for important measurements, yet extra data imply an overdetermined system. We need to develop the mathematics to handle the overdetermined system properly.

The quantity^{11,12}

$$\mathbf{r}(\mathbf{x}) \equiv A\mathbf{x} - \mathbf{b} \quad (13.26)$$

measures how nearly some candidate solution \mathbf{x} solves the system (13.25). We call this quantity the *residual*, and the smaller, the better. More precisely, the smaller the nonnegative real scalar

$$[\mathbf{r}(\mathbf{x})]^*[\mathbf{r}(\mathbf{x})] = \sum_i |r_i(\mathbf{x})|^2 \quad (13.27)$$

is, called the *squared residual norm*, the more favorably we regard the candidate solution \mathbf{x} .

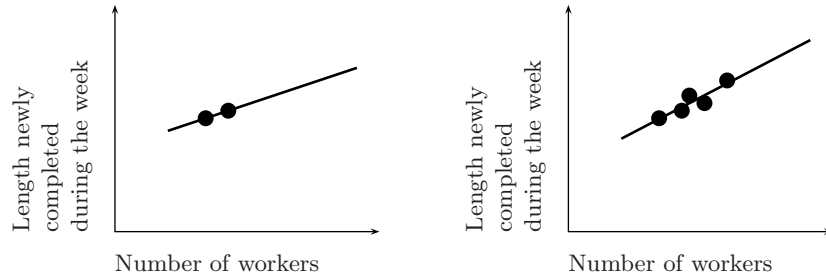
13.6 The Moore-Penrose pseudoinverse and the least-squares problem

A typical problem is to fit a straight line to some data. For example, suppose that we are building-construction contractors with a unionized work force, whose labor union can supply additional, fully trained labor on demand. Suppose further that we are contracted to build a long freeway and have been adding workers to the job in recent weeks to speed construction. On Saturday morning at the end of the second week, we gather and plot the production data on the left of Fig. 13.1. If u_i and b_i respectively represent the number of workers and the length of freeway completed during week i ,

¹¹Alas, the alphabet has only so many letters (see Appendix B). The \mathbf{r} here is unrelated to matrix rank r .

¹²Some books [43] prefer to define $\mathbf{r}(\mathbf{x}) \equiv \mathbf{b} - A\mathbf{x}$, instead.

Figure 13.1: Fitting a line to measured data.



then we can fit a straight line $b = \sigma u + \gamma$ to the measured production data such that

$$\begin{bmatrix} u_1 & 1 \\ u_2 & 1 \end{bmatrix} \begin{bmatrix} \sigma \\ \gamma \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix},$$

inverting the matrix per §§ 13.1 and 13.2 to solve for $\mathbf{x} \equiv [\sigma \ \gamma]^T$, in the hope that the resulting line will predict future production accurately.

That is all mathematically irreproachable. By the fifth Saturday however we shall have gathered more production data, plotted on the figure's right, to which we should like to fit a better line to predict production more accurately. The added data present a problem. Statistically, the added data are welcome, but geometrically we need only two points to specify a line; what are we to do with the other three? The five points together overdetermine the linear system

$$\begin{bmatrix} u_1 & 1 \\ u_2 & 1 \\ u_3 & 1 \\ u_4 & 1 \\ u_5 & 1 \end{bmatrix} \begin{bmatrix} \sigma \\ \gamma \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \end{bmatrix}.$$

There is no way to draw a single straight line $b = \sigma u + \gamma$ exactly through all five, for in placing the line we enjoy only two degrees of freedom.

The proper approach is to draw a line among the data points that misses the points as narrowly as possible. More precisely, the proper approach chooses parameters σ and γ to minimize the squared residual norm

$[\mathbf{r}(\mathbf{x})]^*[\mathbf{r}(\mathbf{x})]$ of § 13.5, given that

$$A = \begin{bmatrix} u_1 & 1 \\ u_2 & 1 \\ u_3 & 1 \\ u_4 & 1 \\ u_5 & 1 \\ \vdots & \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} \sigma \\ \gamma \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \\ \vdots \end{bmatrix}.$$

Such parameters constitute a *least-squares* solution.

The matrix A in the example has two columns, data marching on the left, all ones on the right. This is a typical structure for A , but in general any matrix A with any number of columns of any content might arise (because there were more than two relevant variables or because some data merited heavier weight than others, among many further reasons). Whatever matrix A might arise from whatever source, this section attacks the difficult but important problem of approximating optimally a solution to the general, possibly unsolvable linear system (13.25), $A\mathbf{x} \approx \mathbf{b}$.

13.6.1 Least squares in the real domain

The least-squares problem is simplest when the matrix A enjoys full column rank and no complex numbers are involved. In this case, we seek to minimize the squared residual norm

$$\begin{aligned} [\mathbf{r}(\mathbf{x})]^T[\mathbf{r}(\mathbf{x})] &= (A\mathbf{x} - \mathbf{b})^T(A\mathbf{x} - \mathbf{b}) \\ &= \mathbf{x}^T A^T A\mathbf{x} + \mathbf{b}^T \mathbf{b} - (\mathbf{x}^T A^T \mathbf{b} + \mathbf{b}^T A\mathbf{x}) \\ &= \mathbf{x}^T A^T A\mathbf{x} + \mathbf{b}^T \mathbf{b} - 2\mathbf{x}^T A^T \mathbf{b} \\ &= \mathbf{x}^T A^T (A\mathbf{x} - 2\mathbf{b}) + \mathbf{b}^T \mathbf{b}, \end{aligned}$$

in which the transpose is used interchangeably for the adjoint because all the numbers involved happen to be real. The norm is minimized where

$$\frac{d}{d\mathbf{x}} (\mathbf{r}^T \mathbf{r}) = 0$$

(in which $d/d\mathbf{x}$ is the Jacobian operator of § 11.10). A requirement that

$$\frac{d}{d\mathbf{x}} [\mathbf{x}^T A^T (A\mathbf{x} - 2\mathbf{b}) + \mathbf{b}^T \mathbf{b}] = 0$$

comes of combining the last two equations. Differentiating by the Jacobian product rule (11.79) yields the equation

$$\mathbf{x}^T A^T A + [A^T (A\mathbf{x} - 2\mathbf{b})]^T = 0;$$

or, after transposing the equation, rearranging terms and dividing by 2, the simplified equation

$$A^T A \mathbf{x} = A^T \mathbf{b}.$$

Assuming (as warranted by § 13.6.2, next) that the $n \times n$ square matrix $A^T A$ is invertible, the simplified equation implies the approximate but optimal least-squares solution

$$\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b} \quad (13.28)$$

to the unsolvable linear system (13.25) in the restricted but quite typical case that A and \mathbf{b} are real and A has full column rank.

Equation (13.28) plots the line on Fig. 13.1's right. As the reader can see, the line does not pass through all the points, for no line can; but it does pass pretty convincingly among them.

13.6.2 The invertibility of $A^* A$

Section 13.6.1 has assumed correctly but unwarrantedly that the product $A^T A$ were invertible for real A of full column rank. For real A , it happens that $A^T = A^*$, so it only broadens the same assumption to suppose that the product $A^* A$ were invertible for complex A of full column rank.¹³ This subsection warrants the latter assumption, thereby incidentally also warranting the former.

Let A be a complex, $m \times n$ matrix of full column rank $r = n \leq m$. Suppose falsely that $A^* A$ were not invertible but singular. Since the product $A^* A$ is a square, $n \times n$ matrix, this is to suppose (§ 13.1) that the product's rank $r' < n$ were less than full, implying (§ 12.5.4) that its columns (as its rows) depended on one another. This would mean that there existed a nonzero, n -element vector \mathbf{u} for which

$$A^* A \mathbf{u} = 0, \quad I_n \mathbf{u} \neq 0.$$

Left-multiplying by \mathbf{u}^* would give that

$$\mathbf{u}^* A^* A \mathbf{u} = 0, \quad I_n \mathbf{u} \neq 0,$$

or in other words that

$$(A\mathbf{u})^* (A\mathbf{u}) = \sum_{i=1}^n |[A\mathbf{u}]_i|^2 = 0, \quad I_n \mathbf{u} \neq 0.$$

¹³Notice that if A is tall, then $A^* A$ is a compact, $n \times n$ square, whereas AA^* is a big, $m \times m$ square. It is the compact square that concerns this section. The big square is not very interesting and in any case is not invertible.

But this could only be so if

$$A\mathbf{u} = 0, \quad I_n\mathbf{u} \neq 0,$$

impossible when the columns of A are independent. The contradiction proves false the assumption which gave rise to it. The false assumption: that A^*A were singular.

Thus, *the $n \times n$ product A^*A is invertible for any tall or square, $m \times n$ matrix A of full column rank $r = n \leq m$.*

13.6.3 Positive definiteness

An $n \times n$ matrix C is *positive definite* if and only if

$$\Im(\mathbf{u}^*C\mathbf{u}) = 0 \text{ and } \Re(\mathbf{u}^*C\mathbf{u}) > 0 \text{ for all } I_n\mathbf{u} \neq 0. \quad (13.29)$$

As in § 13.6.2, here also when a matrix A has full column rank $r = n \leq m$ the product $\mathbf{u}^*A^*A\mathbf{u} = (A\mathbf{u})^*(A\mathbf{u})$ is real and positive for all nonzero, n -element vectors \mathbf{u} . Thus per (13.29) *the product A^*A is positive definite for any matrix A of full column rank.*

An $n \times n$ matrix C is *nonnegative definite* if and only if

$$\Im(\mathbf{u}^*C\mathbf{u}) = 0 \text{ and } \Re(\mathbf{u}^*C\mathbf{u}) \geq 0 \text{ for all } \mathbf{u}. \quad (13.30)$$

By reasoning like the last paragraph's, *the product A^*A is nonnegative definite for any matrix A whatsoever.*

13.6.4 The Moore-Penrose pseudoinverse

Not every $m \times n$ matrix A enjoys full rank. According to (12.16), however, every $m \times n$ matrix A of rank r can be factored into a product¹⁴

$$A = BC$$

of an $m \times r$ tall or square matrix B and an $r \times n$ broad or square matrix C , both of which factors themselves enjoy full rank r . (If A happens to have full row or column rank, then one can just choose $B = I_m$ or $C = I_n$; but even if A lacks full rank, the Gauss-Jordan decomposition of eqn. 12.2 finds at least the full-rank factorization $B = G_{>}I_r$, $C = I_rG_{<}$.) This being

¹⁴This subsection uses the symbols B and \mathbf{b} for unrelated purposes, which is unfortunate but conventional. See footnote 11.

so, a conjecture seems warranted. Suppose that, inspired by (13.28), we manipulated (13.25) by the successive steps

$$\begin{aligned} A\mathbf{x} &\approx \mathbf{b}, \\ BC\mathbf{x} &\approx \mathbf{b}, \\ (B^*B)^{-1}B^*BC\mathbf{x} &\approx (B^*B)^{-1}B^*\mathbf{b}, \\ C\mathbf{x} &\approx (B^*B)^{-1}B^*\mathbf{b}. \end{aligned}$$

Then suppose that we changed

$$C^*\mathbf{u} \leftarrow \mathbf{x},$$

thus restricting \mathbf{x} to the space addressed by the independent columns of C^* . Continuing,

$$\begin{aligned} CC^*\mathbf{u} &\approx (B^*B)^{-1}B^*\mathbf{b}, \\ \mathbf{u} &\approx (CC^*)^{-1}(B^*B)^{-1}B^*\mathbf{b}. \end{aligned}$$

Changing the variable back and (because we are conjecturing and can do as we like), altering the “ \approx ” sign to “ $=$,”

$$\mathbf{x} = C^*(CC^*)^{-1}(B^*B)^{-1}B^*\mathbf{b}. \quad (13.31)$$

Equation (13.31) has a pleasingly symmetrical form, and we know from § 13.6.2 at least that the two matrices it tries to invert are invertible. So here is our conjecture:

- no \mathbf{x} enjoys a smaller squared residual norm $\mathbf{r}^*\mathbf{r}$ than the \mathbf{x} of (13.31) does; and
- among all \mathbf{x} that enjoy the same, minimal squared residual norm, the \mathbf{x} of (13.31) is strictly least in magnitude.

The conjecture is bold, but if you think about it in the right way it is not unwarranted under the circumstance. After all, (13.31) does resemble (13.28), the latter of which admittedly requires real A of full column rank but does minimize the residual when its requirements are met; and, even if there were more than one \mathbf{x} which minimized the residual, one of them might be smaller than the others: why not the \mathbf{x} of (13.31)? One can but investigate.

The first point of the conjecture is symbolized

$$\mathbf{r}(\mathbf{x}) \leq \mathbf{r}(\mathbf{x} + \Delta\mathbf{x}),$$

where $\Delta \mathbf{x}$ represents the deviation, whether small, moderate or large, of some alternate \mathbf{x} from the \mathbf{x} of (13.31). According to (13.26), this is

$$[A\mathbf{x} - \mathbf{b}]^*[A\mathbf{x} - \mathbf{b}] \leq [(A)(\mathbf{x} + \Delta \mathbf{x}) - \mathbf{b}]^*[(A)(\mathbf{x} + \Delta \mathbf{x}) - \mathbf{b}].$$

Reorganizing,

$$[A\mathbf{x} - \mathbf{b}]^*[A\mathbf{x} - \mathbf{b}] \leq [(A\mathbf{x} - \mathbf{b}) + A\Delta \mathbf{x}]^*[(A\mathbf{x} - \mathbf{b}) + A\Delta \mathbf{x}].$$

Distributing factors and canceling like terms,

$$0 \leq (A\mathbf{x} - \mathbf{b})^*A\Delta \mathbf{x} + \Delta \mathbf{x}^*A^*(A\mathbf{x} - \mathbf{b}) + \Delta \mathbf{x}^*A^*A\Delta \mathbf{x}.$$

But according to (13.31) and the full-rank factorization $A = BC$,

$$\begin{aligned} A^*(A\mathbf{x} - \mathbf{b}) &= A^*A\mathbf{x} - A^*\mathbf{b} \\ &= [C^*B^*][BC][C^*(CC^*)^{-1}(B^*B)^{-1}B^*\mathbf{b}] - [C^*B^*][\mathbf{b}] \\ &= C^*(B^*B)(CC^*)^{-1}(B^*B)^{-1}B^*\mathbf{b} - C^*B^*\mathbf{b} \\ &= C^*B^*\mathbf{b} - C^*B^*\mathbf{b} = 0, \end{aligned}$$

which reveals two of the inequality's remaining three terms to be zero, leaving an assertion that

$$0 \leq \Delta \mathbf{x}^*A^*A\Delta \mathbf{x}.$$

Each step in the present paragraph is reversible,¹⁵ so the assertion in the last form is logically equivalent to the conjecture's first point, with which the paragraph began. Moreover, the assertion in the last form is correct because the product of any matrix and its adjoint according to § 13.6.3 is a nonnegative definite operator, thus establishing the conjecture's first point.

The conjecture's first point, now established, has it that no $\mathbf{x} + \Delta \mathbf{x}$ enjoys a smaller squared residual norm than the \mathbf{x} of (13.31) does. It does not claim that no $\mathbf{x} + \Delta \mathbf{x}$ enjoys the same, minimal squared residual norm. The latter case is symbolized

$$\mathbf{r}(\mathbf{x}) = \mathbf{r}(\mathbf{x} + \Delta \mathbf{x});$$

or equivalently by the last paragraph's logic,

$$0 = \Delta \mathbf{x}^*A^*A\Delta \mathbf{x};$$

or in other words,

$$A\Delta \mathbf{x} = 0.$$

¹⁵The paragraph might inscrutably but logically instead have ordered the steps in reverse as in §§ 6.3.2 and 9.5. See Ch. 6's footnote 14.

But $A = BC$, so this is to claim that

$$B(C \Delta \mathbf{x}) = 0,$$

which since B has full column rank is possible only if

$$C \Delta \mathbf{x} = 0.$$

Considering the product $\Delta \mathbf{x}^* \mathbf{x}$ in light of (13.31) and the last equation, we observe that

$$\begin{aligned} \Delta \mathbf{x}^* \mathbf{x} &= \Delta \mathbf{x}^* [C^* (CC^*)^{-1} (B^* B)^{-1} B^* \mathbf{b}] \\ &= [C \Delta \mathbf{x}]^* [(CC^*)^{-1} (B^* B)^{-1} B^* \mathbf{b}], \end{aligned}$$

which is to observe that

$$\Delta \mathbf{x}^* \mathbf{x} = 0$$

for any $\Delta \mathbf{x}$ for which $\mathbf{x} + \Delta \mathbf{x}$ achieves minimal squared residual norm.

Returning attention to the conjecture, its second point is symbolized

$$\mathbf{x}^* \mathbf{x} < (\mathbf{x} + \Delta \mathbf{x})^* (\mathbf{x} + \Delta \mathbf{x})$$

for any

$$\Delta \mathbf{x} \neq 0$$

for which $\mathbf{x} + \Delta \mathbf{x}$ achieves minimal squared residual norm (note that it's “<” this time, not “≤” as in the conjecture's first point). Distributing factors and canceling like terms,

$$0 < \mathbf{x}^* \Delta \mathbf{x} + \Delta \mathbf{x}^* \mathbf{x} + \Delta \mathbf{x}^* \Delta \mathbf{x}.$$

But the last paragraph has found that $\Delta \mathbf{x}^* \mathbf{x} = 0$ for precisely such $\Delta \mathbf{x}$ as we are considering here, so the last inequality reduces to read

$$0 < \Delta \mathbf{x}^* \Delta \mathbf{x},$$

which naturally for $\Delta \mathbf{x} \neq 0$ is true. Since each step in the paragraph is reversible, reverse logic establishes the conjecture's second point.

With both its points established, the conjecture is true.

If $A = BC$ is a full-rank factorization, then the matrix¹⁶

$$A^\dagger \equiv C^* (CC^*)^{-1} (B^* B)^{-1} B^* \quad (13.32)$$

¹⁶Some books print A^\dagger as A^+ .

of (13.31) is called the *Moore-Penrose pseudoinverse* of A , more briefly the *pseudoinverse* of A . Whether underdetermined, exactly determined, overdetermined or even degenerate, every matrix has a Moore-Penrose pseudoinverse. Yielding the optimal approximation

$$\mathbf{x} = A^\dagger \mathbf{b}, \quad (13.33)$$

the Moore-Penrose solves the linear system (13.25) as well as the system can be solved—exactly if possible, with minimal squared residual norm if impossible. If A is square and invertible, then the Moore-Penrose $A^\dagger = A^{-1}$ is just the inverse, and then of course (13.33) solves the system uniquely and exactly. Nothing can solve the system uniquely if A has broad shape but the Moore-Penrose still solves the system exactly in that case as long as A has full row rank, moreover minimizing the solution's squared magnitude $\mathbf{x}^* \mathbf{x}$ (which the solution of eqn. 13.23 fails to do). If A lacks full row rank, then the Moore-Penrose solves the system as nearly as the system can be solved (as in Fig. 13.1) and as a side-benefit also minimizes $\mathbf{x}^* \mathbf{x}$. The Moore-Penrose is thus a general-purpose solver and approximator for linear systems. It is a significant discovery.¹⁷

13.7 The multivariate Newton-Raphson iteration

When we first met the Newton-Raphson iteration in § 4.8 we lacked the matrix notation and algebra to express and handle vector-valued functions adeptly. Now that we have the notation and algebra we can write down the multivariate Newton-Raphson iteration almost at once.

The iteration approximates the nonlinear vector function $\mathbf{f}(\mathbf{x})$ by its tangent

$$\tilde{\mathbf{f}}_k(\mathbf{x}) = \mathbf{f}(\mathbf{x}_k) + \left[\frac{d}{d\mathbf{x}} \mathbf{f}(\mathbf{x}) \right]_{\mathbf{x}=\mathbf{x}_k} (\mathbf{x} - \mathbf{x}_k),$$

where $d\mathbf{f}/d\mathbf{x}$ is the Jacobian derivative of § 11.10. It then approximates the root \mathbf{x}_{k+1} as the point at which $\tilde{\mathbf{f}}_k(\mathbf{x}_{k+1}) = 0$:

$$\tilde{\mathbf{f}}_k(\mathbf{x}_{k+1}) = 0 = \mathbf{f}(\mathbf{x}_k) + \left[\frac{d}{d\mathbf{x}} \mathbf{f}(\mathbf{x}) \right]_{\mathbf{x}=\mathbf{x}_k} (\mathbf{x}_{k+1} - \mathbf{x}_k).$$

Solving for \mathbf{x}_{k+1} (approximately if necessary), we have that

$$\mathbf{x}_{k+1} = \mathbf{x} - \left[\frac{d}{d\mathbf{x}} \mathbf{f}(\mathbf{x}) \right]_{\mathbf{x}=\mathbf{x}_k}^\dagger \mathbf{f}(\mathbf{x}) \Big|_{\mathbf{x}=\mathbf{x}_k}, \quad (13.34)$$

¹⁷[2, § 3.3][33, “Moore-Penrose generalized inverse”]

where $[\cdot]^\dagger$ is the Moore-Penrose pseudoinverse of § 13.6—which is just the ordinary inverse $[\cdot]^{-1}$ of § 13.1 if \mathbf{f} and \mathbf{x} happen each to have the same number of elements. Refer to § 4.8 and Fig. 4.5.¹⁸

Despite the Moore-Penrose notation of (13.34), the Newton-Raphson iteration is not normally meant to be applied at a value of \mathbf{x} for which the Jacobian is degenerate. The iteration intends rather in light of (13.32) that

$$\left[\frac{d}{d\mathbf{x}} \mathbf{f}(\mathbf{x}) \right]^\dagger = \begin{cases} [d\mathbf{f}/d\mathbf{x}]^* ([d\mathbf{f}/d\mathbf{x}] [d\mathbf{f}/d\mathbf{x}]^*)^{-1} & \text{if } r = m \leq n, \\ [d\mathbf{f}/d\mathbf{x}]^{-1} & \text{if } r = m = n, \\ ([d\mathbf{f}/d\mathbf{x}]^* [d\mathbf{f}/d\mathbf{x}])^{-1} [d\mathbf{f}/d\mathbf{x}]^* & \text{if } r = n \leq m, \end{cases} \quad (13.35)$$

where $B = I_m$ in the first case and $C = I_n$ in the last. It does not intend to use the full (13.32). If both $r < m$ and $r < n$ —which is to say, if the Jacobian is degenerate—then (13.35) fails, as though the curve of Fig. 4.5 ran horizontally at the test point—when one quits, restarting the iteration from another point.

13.8 The dot product

The *dot product* of two vectors, also called the *inner product*,¹⁹ is the product of the two vectors to the extent that they run in the same direction. It is written

$$\mathbf{a} \cdot \mathbf{b}.$$

In general,

$$\mathbf{a} \cdot \mathbf{b} = (a_1 \mathbf{e}_1 + a_2 \mathbf{e}_2 + \cdots + a_n \mathbf{e}_n) \cdot (b_1 \mathbf{e}_1 + b_2 \mathbf{e}_2 + \cdots + b_n \mathbf{e}_n).$$

But if the dot product is to mean anything, it must be that

$$\mathbf{e}_i \cdot \mathbf{e}_j = \delta_{ij}. \quad (13.36)$$

Therefore,

$$\mathbf{a} \cdot \mathbf{b} = a_1 b_1 + a_2 b_2 + \cdots + a_n b_n;$$

¹⁸[35]

¹⁹The term *inner product* is often used to indicate a broader class of products than the one defined here, especially in some of the older literature. Where used, the notation usually resembles $\langle \mathbf{a}, \mathbf{b} \rangle$ or (\mathbf{b}, \mathbf{a}) , both of which mean $\mathbf{a}^* \cdot \mathbf{b}$ (or, more broadly, some similar product), except that which of \mathbf{a} and \mathbf{b} is conjugated depends on the author. Most recently, at least in the author's country, the usage $\langle \mathbf{a}, \mathbf{b} \rangle \equiv \mathbf{a}^* \cdot \mathbf{b}$ seems to be emerging as standard where the dot is not used [2, § 3.1][13, Ch. 4]. This book prefers the dot.

or, more concisely,

$$\mathbf{a} \cdot \mathbf{b} \equiv \mathbf{a}^T \mathbf{b} = \sum_{j=-\infty}^{\infty} a_j b_j. \quad (13.37)$$

The dot notation does not worry whether its arguments are column or row vectors, incidentally:

$$\mathbf{a} \cdot \mathbf{b} = \mathbf{a} \cdot \mathbf{b}^T = \mathbf{a}^T \cdot \mathbf{b} = \mathbf{a}^T \cdot \mathbf{b}^T = \mathbf{a}^T \mathbf{b}.$$

That is, if either vector is wrongly oriented, the notation implicitly reorients it before using it. (The more orderly notation $\mathbf{a}^T \mathbf{b}$ by contrast assumes that both are proper column vectors.)

Where vectors may have complex elements, usually one is not interested in $\mathbf{a} \cdot \mathbf{b}$ so much as in

$$\mathbf{a}^* \cdot \mathbf{b} \equiv \mathbf{a}^* \mathbf{b} = \sum_{j=-\infty}^{\infty} a_j^* b_j. \quad (13.38)$$

The reason is that

$$\Re(\mathbf{a}^* \cdot \mathbf{b}) = \Re(\mathbf{a}) \cdot \Re(\mathbf{b}) + \Im(\mathbf{a}) \cdot \Im(\mathbf{b}),$$

with the product of the imaginary parts added not subtracted, thus honoring the right Argand sense of “the product of the two vectors to the extent that they run in the same direction.”

By the Pythagorean theorem, the dot product

$$|\mathbf{a}|^2 = \mathbf{a}^* \cdot \mathbf{a} \quad (13.39)$$

gives the square of a vector’s magnitude, always real, never negative. The unit vector in \mathbf{a} ’s direction then is

$$\hat{\mathbf{a}} \equiv \frac{\mathbf{a}}{|\mathbf{a}|} = \frac{\mathbf{a}}{\sqrt{\mathbf{a}^* \cdot \mathbf{a}}}, \quad (13.40)$$

from which

$$|\hat{\mathbf{a}}|^2 = \hat{\mathbf{a}}^* \cdot \hat{\mathbf{a}} = 1. \quad (13.41)$$

When two vectors do not run in the same direction at all, such that

$$\mathbf{a}^* \cdot \mathbf{b} = 0, \quad (13.42)$$

the two vectors are said to lie *orthogonal* to one another. Geometrically this puts them at right angles. For other angles θ between two vectors,

$$\hat{\mathbf{a}}^* \cdot \hat{\mathbf{b}} = \cos \theta, \quad (13.43)$$

which formally defines the angle θ even when \mathbf{a} and \mathbf{b} have more than three elements each.

13.9 The orthogonal complement

The $m \times (m - r)$ kernel (§ 13.3)²⁰

$$A^\perp \equiv A^{*K} \quad (13.44)$$

is an interesting matrix. By definition of the kernel, the columns of A^{*K} are the independent vectors \mathbf{u}_j for which $A^*\mathbf{u}_j = 0$, which—inasmuch as the rows of A^* are the adjoints of the *columns* of A —is possible only when each \mathbf{u}_j lies orthogonal to every column of A . This says that the columns of $A^\perp \equiv A^{*K}$ address the complete space of vectors that lie orthogonal to A 's columns, such that

$$A^{\perp*}A = 0 = A^*A^\perp. \quad (13.45)$$

The matrix A^\perp is called the *orthogonal complement*²¹ or *perpendicular matrix* to A .

Among other uses, the orthogonal complement A^\perp supplies the columns A lacks to reach full row rank. Properties include that

$$\begin{aligned} A^{*K} &= A^\perp, \\ A^{*\perp} &= A^K. \end{aligned} \quad (13.46)$$

13.10 Gram-Schmidt orthonormalization

If a vector $\mathbf{x} = A^K \mathbf{a}$ belongs to a kernel space A^K (§ 13.3), then so equally does any $\alpha \mathbf{x}$. If the vectors $\mathbf{x}_1 = A^K \mathbf{a}_1$ and $\mathbf{x}_2 = A^K \mathbf{a}_2$ both belong, then so does $\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2$. If I claim $A^K = [3 \ 4 \ 5; -1 \ 1 \ 0]^T$ to represent a kernel, then you are not mistaken arbitrarily to rescale each column of my A^K by a separate nonzero factor, instead for instance representing the same kernel as $A^K = [6 \ 8 \ 0; \frac{1}{7} \ -\frac{1}{7} \ 0]^T$. Kernel vectors have no inherent scale. Style generally asks the applied mathematician to remove the false appearance of scale by using (13.40) *to normalize* the columns of a kernel matrix to unit magnitude before reporting them. The same goes for the eigenvectors of Ch. 14 to come.

Where a kernel matrix A^K has two or more columns (or a repeated eigenvalue has two or more eigenvectors), style generally asks the applied mathematician not only to normalize but also *to orthogonalize* the columns

²⁰The symbol A^\perp [20][2][27] can be pronounced “A perp,” short for “A perpendicular,” since by (13.45) A^\perp is in some sense perpendicular to A .

If we were really precise, we might write not A^\perp but $A^{\perp(m)}$. Refer to footnote 5.

²¹[20, § 3.VI.3]

before reporting them. One orthogonalizes a vector \mathbf{b} with respect to a vector \mathbf{a} by subtracting from \mathbf{b} a multiple of \mathbf{a} such that

$$\begin{aligned}\mathbf{a}^* \cdot \mathbf{b}_\perp &= 0, \\ \mathbf{b}_\perp &\equiv \mathbf{b} - \beta \mathbf{a},\end{aligned}$$

where the symbol \mathbf{b}_\perp represents the orthogonalized vector. Substituting the second of these equations into the first and solving for β yields

$$\beta = \frac{\mathbf{a}^* \cdot \mathbf{b}}{\mathbf{a}^* \cdot \mathbf{a}}.$$

Hence,

$$\begin{aligned}\mathbf{a}^* \cdot \mathbf{b}_\perp &= 0, \\ \mathbf{b}_\perp &\equiv \mathbf{b} - \frac{\mathbf{a}^* \cdot \mathbf{b}}{\mathbf{a}^* \cdot \mathbf{a}} \mathbf{a}.\end{aligned}\tag{13.47}$$

But according to (13.40), $\mathbf{a} = \hat{\mathbf{a}}\sqrt{\mathbf{a}^* \cdot \mathbf{a}}$; and according to (13.41), $\hat{\mathbf{a}}^* \cdot \hat{\mathbf{a}} = 1$; so,

$$\mathbf{b}_\perp = \mathbf{b} - \hat{\mathbf{a}}(\hat{\mathbf{a}}^* \cdot \mathbf{b});\tag{13.48}$$

or, in matrix notation,

$$\mathbf{b}_\perp = \mathbf{b} - \hat{\mathbf{a}}(\hat{\mathbf{a}}^*)(\mathbf{b}).$$

This is arguably better written

$$\mathbf{b}_\perp = [I - (\hat{\mathbf{a}})(\hat{\mathbf{a}}^*)] \mathbf{b}\tag{13.49}$$

(observe that it's $[\hat{\mathbf{a}}][\hat{\mathbf{a}}^*]$, a matrix, rather than the scalar $[\hat{\mathbf{a}}^*][\hat{\mathbf{a}}]$).

One *orthonormalizes* a set of vectors by orthogonalizing them with respect to one another, then by normalizing each of them to unit magnitude. The procedure to orthonormalize several vectors

$$\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_n\}$$

therefore is as follows. First, normalize \mathbf{x}_1 by (13.40); call the result $\hat{\mathbf{x}}_{1\perp}$. Second, orthogonalize \mathbf{x}_2 with respect to $\hat{\mathbf{x}}_{1\perp}$ by (13.48) or (13.49), then normalize it; call the result $\hat{\mathbf{x}}_{2\perp}$. Third, orthogonalize \mathbf{x}_3 with respect to $\hat{\mathbf{x}}_{1\perp}$ then to $\hat{\mathbf{x}}_{2\perp}$, then normalize it; call the result $\hat{\mathbf{x}}_{3\perp}$. Proceed in this manner through the several \mathbf{x}_j . Symbolically,

$$\begin{aligned}\hat{\mathbf{x}}_{j\perp} &= \frac{\mathbf{x}_{j\perp}}{\sqrt{\mathbf{x}_{j\perp}^* \mathbf{x}_{j\perp}}}, \\ \mathbf{x}_{j\perp} &\equiv \left[\prod_{i=1}^{j-1} (I - \hat{\mathbf{x}}_{i\perp} \hat{\mathbf{x}}_{i\perp}^*) \right] \mathbf{x}_j.\end{aligned}\tag{13.50}$$

By the vector replacement principle of § 12.4 in light of (13.47), the resulting orthonormal set of vectors

$$\{\hat{\mathbf{x}}_{1\perp}, \hat{\mathbf{x}}_{2\perp}, \hat{\mathbf{x}}_{3\perp}, \dots, \hat{\mathbf{x}}_{n\perp}\}$$

addresses the same space as did the original set.

Orthonormalization naturally works equally for any linearly independent set of vectors, not only for kernel vectors or eigenvectors. By the technique, one can conveniently replace a set of independent vectors by an equivalent, neater, orthonormal set which addresses precisely the same space.

13.10.1 Efficient implementation

To turn an equation like the latter line of (13.50) into an efficient numerical algorithm sometimes demands some extra thought, in perspective of whatever it happens to be that one is trying to accomplish. If all one wants is some vectors orthonormalized, then the equation as written is neat but is overkill because the product $\hat{\mathbf{x}}_{i\perp} \hat{\mathbf{x}}_{i\perp}^*$ is a matrix, whereas the product $\hat{\mathbf{x}}_{i\perp}^* \mathbf{x}_j$ implied by (13.48) is just a scalar. Fortunately, one need not apply the latter line of (13.50) exactly as written. One can instead introduce intermediate vectors \mathbf{x}_{ji} , representing the \prod multiplication in the admittedly messier form

$$\begin{aligned} \mathbf{x}_{j1} &\equiv \mathbf{x}_j, \\ \mathbf{x}_{j(i+1)} &\equiv \mathbf{x}_{ji} - (\hat{\mathbf{x}}_{i\perp}^* \cdot \mathbf{x}_{ji}) \hat{\mathbf{x}}_{i\perp}, \\ \mathbf{x}_{j\perp} &= \mathbf{x}_{jj}. \end{aligned} \tag{13.51}$$

Besides obviating the matrix $I - \hat{\mathbf{x}}_{i\perp} \hat{\mathbf{x}}_{i\perp}^*$ and the associated matrix multiplication, the messier form (13.51) has the significant additional practical virtue that it lets one forget each intermediate vector \mathbf{x}_{ji} immediately after using it. (A well written orthonormalizing computer program reserves memory for one intermediate vector only, which memory it repeatedly overwrites—and, actually, probably does not even reserve that much, working rather in the memory space it has already reserved for $\hat{\mathbf{x}}_{j\perp}$.)²²

Other equations one algorithmizes can likewise benefit from thoughtful rendering.

13.10.2 The Gram-Schmidt decomposition

The orthonormalization technique this section has developed is named the *Gram-Schmidt process*. One can turn it into the *Gram-Schmidt decomposi-*

²²[46, “Gram-Schmidt process,” 04:48, 11 August 2007]

tion

$$\begin{aligned} A &= QR = QUDS, \\ R &\equiv UDS, \end{aligned} \tag{13.52}$$

also called the *orthonormalizing* or *QR decomposition*, by an algorithm that somewhat resembles the Gauss-Jordan algorithm of § 12.3.3; except that (12.4) here becomes

$$A = \tilde{Q}\tilde{U}\tilde{D}\tilde{S} \tag{13.53}$$

and initially $\tilde{Q} \leftarrow A$. By elementary column operations based on (13.50) and (13.51), the algorithm gradually transforms \tilde{Q} into a dimension-limited, $m \times r$ matrix Q of orthonormal columns, distributing the inverse elementaries to \tilde{U} , \tilde{D} and \tilde{S} according to Table 12.1—where the latter three working matrices ultimately become the extended-operational factors U , D and S of (13.52).

Borrowing the language of computer science we observe that the indices i and j of (13.50) and (13.51) imply a two-level nested loop, one level looping over j and the other over i . The equations suggest *j-major nesting*, with the loop over j at the outer level and the loop over i at the inner, such that the several (i, j) index pairs occur in the sequence (reading left to right then top to bottom)

$$\begin{array}{cccc} (1, 2) & & & \\ (1, 3) & (2, 3) & & \\ (1, 4) & (2, 4) & (3, 4) & \\ \dots & \dots & \dots & \ddots \end{array}$$

In reality, however, (13.51)'s middle line requires only that no $\hat{\mathbf{x}}_{i\perp}$ be used before it is fully calculated; otherwise that line does not care which (i, j) pair follows which. The *i-major nesting*

$$\begin{array}{cccc} (1, 2) & (1, 3) & (1, 4) & \dots \\ & (2, 3) & (2, 4) & \dots \\ & & (3, 4) & \dots \\ & & & \ddots \end{array}$$

bringing the very same index pairs in a different sequence, is just as valid. We choose *i-major nesting* on the subtle ground that it affords better information to the choice of column index p during the algorithm's step 3.

The algorithm, in detail:

1. Begin by initializing

$$\begin{aligned}\tilde{Q} &\leftarrow A, \\ \tilde{U} &\leftarrow I, \\ \tilde{D} &\leftarrow I, \\ \tilde{S} &\leftarrow I, \\ i &\leftarrow 1.\end{aligned}$$

2. (Besides arriving at this point from step 1 above, the algorithm also reënters here from step 9 below.) Observe that \tilde{U} enjoys the major partial unit triangular form $L^{\{i-1\}T}$ (§ 11.8.5); that \tilde{D} is a general scaling operator (§ 11.7.2) with $\tilde{d}_{jj} = 1$ for all $j \geq i$; that \tilde{S} is permutor (§ 11.7.1); and that the first through $(i-1)$ th columns of \tilde{Q} consist of mutually orthonormal unit vectors.
3. Choose a column $p \geq i$ of \tilde{Q} containing at least one nonzero element. (The simplest choice is perhaps $p = i$ as long as the i th column does not happen to be null, but one might instead prefer to choose the column of greatest magnitude, or to choose randomly, among other heuristics.) If \tilde{Q} is null in and rightward of its i th column such that no column $p \geq i$ remains available to choose, then skip directly to step 10.
4. Observing that (13.53) can be expanded to read

$$\begin{aligned}A &= \left(\tilde{Q}T_{[i \leftrightarrow p]}\right)\left(T_{[i \leftrightarrow p]}\tilde{U}T_{[i \leftrightarrow p]}\right)\left(T_{[i \leftrightarrow p]}\tilde{D}T_{[i \leftrightarrow p]}\right)\left(T_{[i \leftrightarrow p]}\tilde{S}\right) \\ &= \left(\tilde{Q}T_{[i \leftrightarrow p]}\right)\left(T_{[i \leftrightarrow p]}\tilde{U}T_{[i \leftrightarrow p]}\right)\tilde{D}\left(T_{[i \leftrightarrow p]}\tilde{S}\right),\end{aligned}$$

where the latter line has applied a rule from Table 12.1, interchange the chosen p th column to the i th position by

$$\begin{aligned}\tilde{Q} &\leftarrow \tilde{Q}T_{[i \leftrightarrow p]}, \\ \tilde{U} &\leftarrow T_{[i \leftrightarrow p]}\tilde{U}T_{[i \leftrightarrow p]}, \\ \tilde{S} &\leftarrow T_{[i \leftrightarrow p]}\tilde{S}.\end{aligned}$$

5. Observing that (13.53) can be expanded to read

$$A = \left(\tilde{Q}T_{(1/\alpha)[i]}\right)\left(T_{\alpha[i]}\tilde{U}T_{(1/\alpha)[i]}\right)\left(T_{\alpha[i]}\tilde{D}\right)\tilde{S},$$

normalize the i th column of \tilde{Q} by

$$\begin{aligned}\tilde{Q} &\leftarrow \tilde{Q}T_{(1/\alpha)[i]}, \\ \tilde{U} &\leftarrow T_{\alpha[i]}\tilde{U}T_{(1/\alpha)[i]}, \\ \tilde{D} &\leftarrow T_{\alpha[i]}\tilde{D},\end{aligned}$$

where

$$\alpha = \sqrt{[\tilde{Q}]_{*i}^* \cdot [\tilde{Q}]_{*i}}.$$

6. Initialize

$$j \leftarrow i + 1.$$

7. (Besides arriving at this point from step 6 above, the algorithm also reënters here from step 8 below.) If $j > n$ then skip directly to step 9. Otherwise, observing that (13.53) can be expanded to read

$$A = (\tilde{Q}T_{-\beta[ij]})(T_{\beta[ij]}\tilde{U})\tilde{D}\tilde{S},$$

orthogonalize the j th column of \tilde{Q} per (13.51) with respect to the i th column by

$$\begin{aligned}\tilde{Q} &\leftarrow \tilde{Q}T_{-\beta[ij]}, \\ \tilde{U} &\leftarrow T_{\beta[ij]}\tilde{U},\end{aligned}$$

where

$$\beta = [\tilde{Q}]_{*i}^* \cdot [\tilde{Q}]_{*j}.$$

8. Increment

$$j \leftarrow j + 1$$

and return to step 7.

9. Increment

$$i \leftarrow i + 1$$

and return to step 2.

10. Let

$$\begin{aligned}Q &\equiv \tilde{Q}, \\ U &\equiv \tilde{U}, \\ D &\equiv \tilde{D}, \\ S &\equiv \tilde{S}, \\ r &= i - 1.\end{aligned}$$

End.

Though the Gram-Schmidt algorithm broadly resembles the Gauss-Jordan, at least two significant differences stand out: (i) the Gram-Schmidt is one-sided because it operates only on the columns of \tilde{Q} , never on the rows; (ii) since Q is itself dimension-limited, the Gram-Schmidt decomposition (13.52) needs and has no explicit factor I_r .

As in § 12.5.7, here also one sometimes prefers that $S = I$. The algorithm optionally supports this preference if the $m \times n$ matrix A has full column rank $r = n$, when null columns cannot arise, if one always chooses $p = i$ during the algorithm's step 3. Such optional discipline maintains $S = I$ when desired.

Whether $S = I$ or not, the matrix $Q = QI_r$ has only r columns, so one can write (13.52) as

$$A = (QI_r)(R).$$

Reassociating factors, this is

$$A = (Q)(I_r R), \quad (13.54)$$

which per (12.16) is a proper full-rank factorization with which one can compute the pseudoinverse A^\dagger of A (see eqn. 13.32, above; but see also eqn. 13.63, below).

If the Gram-Schmidt decomposition (13.52) looks useful, it is even more useful than it looks. The most interesting of its several factors is the $m \times r$ orthonormalized matrix Q , whose orthonormal columns address the same space the columns of A themselves address. If Q reaches the maximum possible rank $r = m$, achieving square, $m \times m$ shape, then it becomes a *unitary matrix*—the subject of § 13.11. Before treating the unitary matrix, however, let us pause to extract a kernel from the Gram-Schmidt decomposition in § 13.10.3, next.

13.10.3 The Gram-Schmidt kernel formula

Like the Gauss-Jordan decomposition in (13.7), the Gram-Schmidt decomposition too brings a kernel formula. To develop and apply it, one decomposes an $m \times n$ matrix

$$A = QR \quad (13.55)$$

per the Gram-Schmidt (13.52) and its algorithm in § 13.10.2. Observing that the r independent columns of the $m \times r$ matrix Q address the same space the columns of A address, one then constructs the $m \times (r + m)$ matrix

$$A' \equiv Q + I_m H_{-r} = \begin{bmatrix} Q & I_m \end{bmatrix} \quad (13.56)$$

and decomposes it too,

$$A' = Q'R', \quad (13.57)$$

again by Gram-Schmidt—with the difference that, this time, one initializes

$$i \leftarrow r + 1$$

rather than $i \leftarrow 1$ at the algorithm's step 1; on the ground that the earlier Gram-Schmidt application of (13.55) has already orthonormalized first r columns of A' , which columns, after all, are just Q . The resulting $m \times m$, full-rank square matrix

$$Q' = Q + A^\perp H_{-r} = \begin{bmatrix} Q & A^\perp \end{bmatrix} \quad (13.58)$$

consists of

- r columns on the left that address the same space the columns of A address and
- $m-r$ columns on the right that give a complete orthogonal complement (§ 13.9) A^\perp of A .

Each column has unit magnitude and conveniently lies orthogonal to every other column, left and right.

Equation (13.58) is probably the more useful form, but the *Gram-Schmidt kernel formula* as such,

$$A^{*K} = A^\perp = Q'H_r I_{m-r}, \quad (13.59)$$

extracts the rightward columns that express the kernel, not of A , but of A^* . To compute the kernel of a matrix B by Gram-Schmidt one sets $A = B^*$ and applies (13.55) through (13.59). Refer to (13.46).

In either the form (13.58) or the form (13.59), the Gram-Schmidt kernel formula does everything the Gauss-Jordan kernel formula (13.7) does and in at least one sense does it better; for, if one wants a Gauss-Jordan kernel orthonormalized, then one must orthonormalize it as an extra step, whereas the Gram-Schmidt kernel comes already orthonormalized.

Being square, the $m \times m$ matrix Q' is a unitary matrix, as the last paragraph of § 13.10.2 has alluded. The unitary matrix is the subject of § 13.11 that follows.

13.11 The unitary matrix

When the orthonormalized matrix Q of the Gram-Schmidt decomposition (13.52) is square, having the maximum possible rank $r = m$, it brings one property so interesting that the property merits a section of its own. The property is that

$$Q^*Q = I_m = QQ^*. \quad (13.60)$$

The reason that $Q^*Q = I_m$ is that Q 's columns are orthonormal, and that the very definition of orthonormality demands that the dot product $[Q]_{*i}^* \cdot [Q]_{*j}$ of orthonormal columns be zero unless $i = j$, when the dot product of a unit vector with itself is unity. That $I_m = QQ^*$ is unexpected, however, until one realizes²³ that the equation $Q^*Q = I_m$ characterizes Q^* to be the rank- m inverse of Q , and that § 13.1 lets any rank- m inverse (orthonormal or otherwise) attack just as well from the right as from the left. Thus,

$$Q^{-1} = Q^*, \quad (13.61)$$

a very useful property. A matrix Q that satisfies (13.60), whether derived from the Gram-Schmidt or from elsewhere, is called a *unitary matrix*. (Note that the permutor of § 11.7.1 enjoys the property of eqn. 13.61 precisely because it is unitary.)

One immediate consequence of (13.60) is that *a square matrix with either orthonormal columns or orthonormal rows is unitary and has both*.

The product of two or more unitary matrices is itself unitary if the matrices are of the same dimensionality. To prove it, consider the product

$$Q = Q_a Q_b \quad (13.62)$$

of $m \times m$ unitary matrices Q_a and Q_b . Let the symbols \mathbf{q}_j , \mathbf{q}_{aj} and \mathbf{q}_{bj} respectively represent the j th columns of Q , Q_a and Q_b , and let the symbol q_{bij} represent the i th element of \mathbf{q}_{bj} . By the columnwise interpretation (§ 11.1.3) of matrix multiplication,

$$\mathbf{q}_j = \sum_i q_{bij} \mathbf{q}_{ai}.$$

The adjoint dot product of any two of Q 's columns then is

$$\mathbf{q}_{j'}^* \cdot \mathbf{q}_j = \sum_{i,i'} q_{bi'j'}^* q_{bij} \mathbf{q}_{ai'}^* \cdot \mathbf{q}_{ai}.$$

²³[13, § 4.4]

But $\mathbf{q}_{ai'}^* \cdot \mathbf{q}_{ai} = \delta_{i'i}$ because Q_a is unitary,²⁴ so

$$\mathbf{q}_{j'}^* \cdot \mathbf{q}_j = \sum_i q_{bi'j}^* q_{bij} = \mathbf{q}_{bj'}^* \cdot \mathbf{q}_{bj} = \delta_{j'j},$$

which says neither more nor less than that the columns of Q are orthonormal, which is to say that Q is unitary, as was to be demonstrated.

Unitary operations preserve length. That is, operating on an m -element vector by an $m \times m$ unitary matrix does not alter the vector's magnitude. To prove it, consider the system

$$Q\mathbf{x} = \mathbf{b}.$$

Multiplying the system by its own adjoint yields

$$\mathbf{x}^* Q^* Q \mathbf{x} = \mathbf{b}^* \mathbf{b}.$$

But according to (13.60), $Q^* Q = I_m$; so,

$$\mathbf{x}^* \mathbf{x} = \mathbf{b}^* \mathbf{b},$$

as was to be demonstrated.

Equation (13.61) lets one use the Gram-Schmidt decomposition (13.52) to invert a square matrix as

$$A^{-1} = R^{-1} Q^* = S^* D^{-1} U^{-1} Q^*. \quad (13.63)$$

Unitary extended operators are certainly possible, for if Q is an $m \times m$ dimension-limited matrix, then the extended operator

$$Q_\infty = Q + (I - I_m),$$

which is just Q with ones running out the main diagonal from its active region, itself meets the unitary criterion (13.60) for $m = \infty$.

Unitary matrices are so easy to handle that they can sometimes justify significant effort to convert a model to work in terms of them if possible. We shall meet the unitary matrix again in §§ 14.10 and 14.12.

The chapter as a whole has demonstrated at least in theory (and usually in practice) techniques to solve any linear system characterized by a matrix of finite dimensionality, whatever the matrix's rank or shape. It has explained how to orthonormalize a set of vectors and has derived from the

²⁴This is true only for $1 \leq i \leq m$, but you knew that already.

explanation the useful Gram-Schmidt decomposition. As the chapter's introduction had promised, the matrix has shown its worth here; for without the matrix's notation, arithmetic and algebra most of the chapter's findings would have lain beyond practical reach. And even so, the single most interesting agent of matrix arithmetic remains yet to be treated. This last is the eigenvalue, and it is the subject of Ch. 14, next.

Chapter 14

The eigenvalue

The *eigenvalue* is a scalar by which a square matrix scales a vector without otherwise changing it, such that

$$A\mathbf{v} = \lambda\mathbf{v}.$$

This chapter analyzes the eigenvalue and the associated *eigenvector* it scales.

Before treating the eigenvalue proper, the chapter gathers together from across Chs. 11 through 14 several properties all invertible square matrices share, assembling the properties in § 14.2 and summarizing them there for reference. One of these properties regards the square matrix's *determinant*, which the book has not yet introduced but which the development of the eigenvalue needs. The determinant opens the chapter.

14.1 The determinant

Through Chs. 11, 12 and 13 the theory of the matrix has developed slowly but pretty straightforwardly. Here comes the first unexpected turn.

It begins with an arbitrary-seeming definition. The *determinant* of an $n \times n$ square matrix A is the sum of $n!$ terms, each term the product of n elements, no two elements from the same row or column, terms of positive parity adding to and terms of negative parity subtracting from the sum—a term's parity (§ 11.6) being the parity of the permutor (§ 11.7.1) marking the positions of the term's elements.

Unless you already know about determinants, the definition alone might seem hard to parse, so try this. The inverse of the general 2×2 square matrix

$$A_2 = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix},$$

by the Gauss-Jordan method or any other convenient technique, is found to be

$$A_2^{-1} = \frac{\begin{bmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{bmatrix}}{a_{11}a_{22} - a_{12}a_{21}}.$$

The quantity¹

$$\det A_2 = a_{11}a_{22} - a_{12}a_{21}$$

in the denominator is defined to be the *determinant* of A_2 . Each of the determinant's terms includes one element from each column of the matrix and one from each row, with parity giving the term its \pm sign. The determinant of the general 3×3 square matrix by the same rule is

$$\begin{aligned} \det A_3 = & (a_{11}a_{22}a_{33} + a_{12}a_{23}a_{31} + a_{13}a_{21}a_{32}) \\ & - (a_{13}a_{22}a_{31} + a_{12}a_{21}a_{33} + a_{11}a_{23}a_{32}); \end{aligned}$$

and indeed if we tediously invert such a matrix symbolically, we do find that quantity in the denominator there.

The parity rule merits a more careful description. The parity of a term like $a_{12}a_{23}a_{31}$ is positive because the parity of the permutor, or interchange quasialementary (§ 11.7.1),

$$P = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

marking the positions of the term's elements is positive. The parity of a term like $a_{13}a_{22}a_{31}$ is negative for the same reason. The determinant comprehends all possible such terms, $n!$ in number, half of positive parity and half of negative. (How do we know that exactly half are of positive and half, negative? Answer: by pairing the terms. For every term like $a_{12}a_{23}a_{31}$ whose marking permutor is P , there is a corresponding $a_{13}a_{22}a_{31}$ whose marking permutor is $T_{[1 \leftrightarrow 2]}P$, necessarily of opposite parity. The sole exception to the rule is the 1×1 square matrix, which has no second term to pair.)

¹The determinant $\det A$ used to be written $|A|$, an appropriately terse notation for which this writer confesses some nostalgia. The older notation $|A|$ however unluckily suggests "the magnitude of A ," which though not quite the wrong idea is not quite the right idea, either. The magnitude $|z|$ of a scalar or $|\mathbf{u}|$ of a vector is a real-valued, nonnegative, nonanalytic function of the elements of the quantity in question, whereas the determinant $\det A$ is a complex-valued, analytic function. The book follows convention by denoting the determinant as $\det A$ for this reason among others.

Normally the context implies a determinant's rank n , but the nonstandard notation

$$\det^{(n)} A$$

is available especially to call the rank out, stating explicitly that the determinant has exactly $n!$ terms. (See also §§ 11.3.5 and 11.5 and eqn. 11.49.²)

It is admitted³ that we have not, as yet, actually shown the determinant to be a generally useful quantity; we have merely motivated and defined it. Historically the determinant probably emerged not from abstract considerations but for the mundane reason that the quantity it represents occurred frequently in practice (as in the A_2^{-1} example above). Nothing however logically prevents one from simply defining some quantity which, at first, one merely suspects will later prove useful. So we do here.⁴

14.1.1 Basic properties

The determinant $\det A$ enjoys several useful basic properties.

- If

$$c_{i*} = \begin{cases} a_{i''*} & \text{when } i = i', \\ a_{i'*} & \text{when } i = i'', \\ a_{i*} & \text{otherwise,} \end{cases}$$

or if

$$c_{*j} = \begin{cases} a_{*j''} & \text{when } j = j', \\ a_{*j'} & \text{when } j = j'', \\ a_{*j} & \text{otherwise,} \end{cases}$$

where $i'' \neq i'$ and $j'' \neq j'$, then

$$\det C = -\det A. \quad (14.1)$$

Interchanging rows or columns negates the determinant.

- If

$$c_{i*} = \begin{cases} \alpha a_{i*} & \text{when } i = i', \\ a_{i*} & \text{otherwise,} \end{cases}$$

²And further Ch. 13's footnotes 5 and 20.

³[13, § 1.2]

⁴[13, Ch. 1]

or if

$$c_{*j} = \begin{cases} \alpha a_{*j} & \text{when } j = j', \\ a_{*j} & \text{otherwise,} \end{cases}$$

then

$$\det C = \alpha \det A. \quad (14.2)$$

Scaling a single row or column of a matrix scales the matrix's determinant by the same factor. (Equation 14.2 tracks the linear scaling property of § 7.3.3 and of eqn. 11.2.)

• If

$$c_{i*} = \begin{cases} a_{i*} + b_{i*} & \text{when } i = i', \\ a_{i*} = b_{i*} & \text{otherwise,} \end{cases}$$

or if

$$c_{*j} = \begin{cases} a_{*j} + b_{*j} & \text{when } j = j', \\ a_{*j} = b_{*j} & \text{otherwise,} \end{cases}$$

then

$$\det C = \det A + \det B. \quad (14.3)$$

If one row or column of a matrix C is the sum of the corresponding rows or columns of two other matrices A and B , while the three matrices remain otherwise identical, then the determinant of the one matrix is the sum of the determinants of the other two. (Equation 14.3 tracks the linear superposition property of § 7.3.3 and of eqn. 11.2.)

• If

$$c_{i'*} = 0,$$

or if

$$c_{*j'} = 0,$$

then

$$\det C = 0. \quad (14.4)$$

A matrix with a null row or column also has a null determinant.

• If

$$c_{i''*} = \gamma c_{i'*},$$

or if

$$c_{*j''} = \gamma c_{*j'},$$

where $i'' \neq i'$ and $j'' \neq j'$, then

$$\det C = 0. \quad (14.5)$$

The determinant is zero if one row or column of the matrix is a multiple of another.

- The determinant of the adjoint is just the determinant's conjugate, and the determinant of the transpose is just the determinant itself:

$$\begin{aligned} \det C^* &= (\det C)^*; \\ \det C^T &= \det C. \end{aligned} \quad (14.6)$$

These basic properties are all fairly easy to see if the definition of the determinant is clearly understood. Equations (14.2), (14.3) and (14.4) come because each of the $n!$ terms in the determinant's expansion has exactly one element from row i' or column j' . Equation (14.1) comes because a row or column interchange reverses parity. Equation (14.6) comes because according to § 11.7.1, the permutors P and P^* always have the same parity, and because the adjoint operation individually conjugates each element of C . Finally, (14.5) comes because, in this case, every term in the determinant's expansion finds an equal term of opposite parity to offset it. Or, more formally, (14.5) comes because the following procedure does not alter the matrix: (i) scale row i'' or column j'' by $1/\gamma$; (ii) scale row i' or column j' by γ ; (iii) interchange rows $i' \leftrightarrow i''$ or columns $j' \leftrightarrow j''$. Not altering the matrix, the procedure does not alter the determinant either; and indeed according to (14.2), step (ii)'s effect on the determinant cancels that of step (i). However, according to (14.1), step (iii) negates the determinant. Hence the net effect of the procedure is to negate the determinant—to negate the very determinant the procedure is not permitted to alter. The apparent contradiction can be reconciled only if the determinant is zero to begin with.

From the foregoing properties the following further property can be deduced.

- If

$$c_{i*} = \begin{cases} a_{i*} + \alpha a_{i'*} & \text{when } i = i'', \\ a_{i*} & \text{otherwise,} \end{cases}$$

or if

$$c_{*j} = \begin{cases} a_{*j} + \alpha a_{*j'} & \text{when } j = j'', \\ a_{*j} & \text{otherwise,} \end{cases}$$

where $i'' \neq i'$ and $j'' \neq j'$, then

$$\det C = \det A. \quad (14.7)$$

Adding to a row or column of a matrix a multiple of another row or column does not change the matrix's determinant.

To derive (14.7) for rows (the column proof is similar), one defines a matrix B such that

$$b_{i*} \equiv \begin{cases} \alpha a_{i'*} & \text{when } i = i'', \\ a_{i*} & \text{otherwise.} \end{cases}$$

From this definition, $b_{i''*} = \alpha a_{i'*}$ whereas $b_{i'*} = a_{i'*}$, so

$$b_{i''*} = \alpha b_{i'*},$$

which by (14.5) guarantees that

$$\det B = 0.$$

On the other hand, the three matrices A , B and C differ only in the (i'') th row, where $[C]_{i''*} = [A]_{i''*} + [B]_{i''*}$; so, according to (14.3),

$$\det C = \det A + \det B.$$

Equation (14.7) results from combining the last two equations.

14.1.2 The determinant and the elementary operator

Section 14.1.1 has it that interchanging, scaling or adding rows or columns of a matrix respectively negates, scales or does not alter the matrix's determinant. But the three operations named are precisely the operations of the three elementaries of § 11.4. Therefore,

$$\begin{aligned} \det T_{[i \leftrightarrow j]} A &= -\det A = \det AT_{[i \leftrightarrow j]}, \\ \det T_{\alpha[i]} A &= \alpha \det A = \det AT_{\alpha[j]}, \\ \det T_{\alpha[ij]} A &= \det A = \det AT_{\alpha[ij]}, \\ 1 \leq (i, j) \leq n, \quad i &\neq j, \end{aligned} \quad (14.8)$$

for any $n \times n$ square matrix A . Obviously also,

$$\begin{aligned} \det IA &= \det A = \det AI, \\ \det I_n A &= \det A = \det AI_n, \\ \det I &= 1 = \det I_n. \end{aligned} \quad (14.9)$$

If A is taken to represent an arbitrary product of identity matrices (I_n and/or I) and elementary operators, then a significant consequence of (14.8) and (14.9), applied recursively, is that the determinant of a product is the product of the determinants, at least where identity matrices and elementary operators are concerned. In symbols,⁵

$$\det \left(\prod_k M_k \right) = \prod_k \det M_k, \quad (14.10)$$

$$M_k \in \{I_n, I, T_{[i \leftrightarrow j]}, T_{\alpha[i]}, T_{\alpha[ij]}\},$$

$$1 \leq (i, j) \leq n.$$

This matters because, as the Gauss-Jordan decomposition of § 12.3 has shown, one can build up any square matrix of full rank by applying elementary operators to I_n . Section 14.1.4 will put the rule (14.10) to good use.

14.1.3 The determinant of a singular matrix

Equation (14.8) gives elementary operators the power to alter a matrix's determinant almost arbitrarily—almost arbitrarily, but not quite. What an $n \times n$ elementary operator⁶ cannot do is to change an $n \times n$ matrix's determinant to or from zero. Once zero, a determinant remains zero under the action of elementary operators. Once nonzero, always nonzero. Elementary operators being reversible have no power to breach this barrier.

Another thing $n \times n$ elementaries cannot do according to § 12.5.3 is to change an $n \times n$ matrix's rank. Nevertheless, such elementaries can reduce any $n \times n$ matrix reversibly to I_r , where $r \leq n$ is the matrix's rank, by the Gauss-Jordan algorithm of § 12.3. Equation (14.4) has that the $n \times n$ determinant of I_r is zero if $r < n$, so it follows that the $n \times n$ determinant of every rank- r matrix is similarly zero if $r < n$; and complementarily that the $n \times n$ determinant of a rank- n matrix is never zero. Singular matrices always have zero determinants; full-rank square matrices never do. One can evidently tell the singularity or invertibility of a square matrix from its determinant alone.

⁵Notation like “ \in ” can be too fancy for applied mathematics, but it does help here. The notation $M_k \in \{\dots\}$ restricts M_k to be any of the things between the braces. As it happens though, in this case, (14.11) below is going to erase the restriction.

⁶That is, an elementary operator which honors an $n \times n$ active region. See § 11.3.2.

14.1.4 The determinant of a matrix product

Sections 14.1.2 and 14.1.3 suggest the useful rule that

$$\det AB = \det A \det B. \quad (14.11)$$

To prove the rule, we consider three distinct cases.

The first case is that A is singular. In this case, B acts as a column operator on A , whereas according to § 12.5.2 no operator has the power to promote A in rank. Hence the product AB is no higher in rank than A , which says that AB is no less singular than A , which implies that AB like A has a null determinant. Evidently (14.11) holds in the first case.

The second case is that B is singular. The proof here resembles that of the first case.

The third case is that neither matrix is singular. Here, we use Gauss-Jordan to decompose both matrices into sequences of elementary operators and rank- n identity matrices, for which

$$\begin{aligned} \det AB &= \det \{[A][B]\} \\ &= \det \left\{ \left[\left(\prod T \right) I_n \left(\prod T \right) \right] \left[\left(\prod T \right) I_n \left(\prod T \right) \right] \right\} \\ &= \left(\prod \det T \right) \det I_n \left(\prod \det T \right) \left(\prod \det T \right) \det I_n \left(\prod \det T \right) \\ &= \det \left[\left(\prod T \right) I_n \left(\prod T \right) \right] \det \left[\left(\prod T \right) I_n \left(\prod T \right) \right] \\ &= \det A \det B, \end{aligned}$$

which is a schematic way of pointing out in light of (14.10) merely that since A and B are products of identity matrices and elementaries, the determinant of the product is the product of the determinants.

So it is that (14.11) holds in all three cases, as was to be demonstrated. *The determinant of a matrix product is the product of the matrix determinants.*

14.1.5 Determinants of inverse and unitary matrices

From (14.11) it follows that

$$\det A^{-1} = \frac{1}{\det A} \quad (14.12)$$

because $A^{-1}A = I_n$ and $\det I_n = 1$.

From (14.6) it follows that if Q is a unitary matrix (§ 13.11), then

$$\begin{aligned}\det Q^* \det Q &= 1, \\ |\det Q| &= 1.\end{aligned}\tag{14.13}$$

This reason is that $|\det Q|^2 = (\det Q)^*(\det Q) = \det Q^* \det Q = \det Q^* Q = \det Q^{-1} Q = \det I_n = 1$.

14.1.6 Inverting the square matrix by determinant

The Gauss-Jordan algorithm comfortably inverts concrete matrices of moderate size, but swamps one in nearly interminable algebra when *symbolically* inverting general matrices larger than the A_2 at the section's head. Slogging through the algebra to invert A_3 symbolically nevertheless (the reader need not actually do this unless he desires a long exercise), one quite incidentally discovers a clever way to factor the determinant:

$$\begin{aligned}C^T A &= (\det A) I_n = AC^T; \\ c_{ij} &\equiv \det R_{ij}; \\ [R_{ij}]_{i'j'} &\equiv \begin{cases} 1 & \text{if } i' = i \text{ and } j' = j, \\ 0 & \text{if } i' = i \text{ or } j' = j \text{ but not both,} \\ a_{i'j'} & \text{otherwise.} \end{cases}\end{aligned}\tag{14.14}$$

Pictorially,

$$R_{ij} = \begin{bmatrix} & \vdots & \vdots & \vdots & \vdots & \vdots & \\ \cdots & * & * & 0 & * & * & \cdots \\ \cdots & * & * & 0 & * & * & \cdots \\ \cdots & 0 & 0 & 1 & 0 & 0 & \cdots \\ \cdots & * & * & 0 & * & * & \cdots \\ \cdots & * & * & 0 & * & * & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \end{bmatrix},$$

same as A except in the i th row and j th column. The matrix C , called the *cofactor* of A , then consists of the determinants of the various R_{ij} .

Another way to write (14.14) is

$$[C^T A]_{ij} = (\det A) \delta_{ij} = [AC^T]_{ij},\tag{14.15}$$

which comprises two cases. In the case that $i = j$,

$$[AC^T]_{ij} = [AC^T]_{ii} = \sum_{\ell} a_{i\ell} c_{i\ell} = \sum_{\ell} a_{i\ell} \det R_{i\ell} = \det A = (\det A) \delta_{ij},$$

wherein the equation $\sum_{\ell} a_{i\ell} \det R_{j\ell} = \det A$ states that $\det A$, being a determinant, consists of several terms, each term including one factor from each row of A , where $a_{i\ell}$ provides the i th row and $R_{j\ell}$ provides the other rows.⁷ In the case that $i \neq j$,

$$[AC^T]_{ij} = \sum_{\ell} a_{i\ell} c_{j\ell} = \sum_{\ell} a_{i\ell} \det R_{j\ell} = 0 = (\det A)(0) = (\det A)\delta_{ij},$$

wherein $\sum_{\ell} a_{i\ell} \det R_{j\ell}$ is the determinant, not of A itself, but rather of A with the j th row replaced by a copy of the i th, which according to (14.5) evaluates to zero. Similar equations can be written for $[C^T A]_{ij}$ in both cases. The two cases together prove (14.15), hence also (14.14).

Dividing (14.14) by $\det A$, we have that⁸

$$\begin{aligned} A^{-1}A &= I_n = AA^{-1}, \\ A^{-1} &= \frac{C^T}{\det A}. \end{aligned} \tag{14.16}$$

Equation (14.16) inverts a matrix by determinant. In practice, it inverts small matrices nicely, through about 4×4 dimensionality (the A_2^{-1} equation at the head of the section is just eqn. 14.16 for $n = 2$). It inverts 5×5 and even 6×6 matrices reasonably, too—especially with the help of a computer to do the arithmetic. Though (14.16) still holds in theory for yet larger matrices, and though symbolically it expresses the inverse of an abstract, $n \times n$ matrix concisely whose entries remain unspecified, for concrete matrices much bigger than 4×4 to 6×6 or so its several determinants begin to grow too great and too many for practical calculation. The Gauss-Jordan technique (or even the Gram-Schmidt technique) is preferred to invert concrete matrices above a certain size for this reason.⁹

14.2 Coincident properties

Chs. 11, 12 and 13, plus this chapter up to the present point, have discovered several coincident properties of the invertible $n \times n$ square matrix. One does

⁷This is a bit subtle, but if you actually write out A_3 and its cofactor C_3 symbolically, trying (14.15) on them, then you will soon see what is meant.

⁸Cramer's rule [13, § 1.6], of which the reader may have heard, results from applying (14.16) to (13.4). However, Cramer's rule is really nothing more than (14.16) in a less pleasing form, so this book does not treat Cramer's rule as such.

⁹For very large matrices, even the Gauss-Jordan grows impractical, due to compound floating-point rounding error and the maybe large but nonetheless limited quantity of available computer memory. Iterative techniques ([chapter not yet written]) serve to invert such matrices approximately.

not feel the full impact of the coincidence when these properties are left scattered across the long chapters; so, let us gather and summarize the properties here. A square, $n \times n$ matrix evidently has either all of the following properties or none of them, never some but not others.

- The matrix is invertible (§ 13.1).
- Its rows are linearly independent (§§ 12.1 and 12.3.4).
- Its columns are linearly independent (§ 12.5.4).
- Its columns address the same space the columns of I_n address, and its rows address the same space the rows of I_n address (§ 12.5.7).
- The Gauss-Jordan algorithm reduces it to I_n (§ 12.3.3). (In this, per § 12.5.3, the choice of pivots does not matter.)
- Decomposing it, the Gram-Schmidt algorithm achieves a fully square, unitary, $n \times n$ factor Q (§ 13.10.2).
- It has full rank $r = n$ (§ 12.5.4).
- The linear system $A\mathbf{x} = \mathbf{b}$ it represents has a unique n -element solution \mathbf{x} , given any specific n -element driving vector \mathbf{b} (§ 13.2).
- The determinant $\det A \neq 0$ (§ 14.1.3).
- None of its eigenvalues is zero (§ 14.3, below).

The square matrix which has one of these properties, has all of them. The square matrix which lacks one, lacks all. Assuming exact arithmetic, a square matrix is either invertible, with all that that implies, or singular; never both. The distinction between invertible and singular matrices is theoretically as absolute as (and indeed is analogous to) the distinction between nonzero and zero scalars.

Whether the distinction is always useful is another matter. Usually the distinction is indeed useful, but a matrix can be *almost* singular just as a scalar can be almost zero. Such a matrix is known, among other ways, by its unexpectedly small determinant. Now it is true: in exact arithmetic, a nonzero determinant, no matter how small, implies a theoretically invertible matrix. Practical matrices however often have entries whose values are imprecisely known; and even when they don't, the computers which invert them tend to do arithmetic imprecisely in floating-point. Matrices which live

on the hazy frontier between invertibility and singularity resemble the infinitesimals of § 4.1.1. They are called *ill conditioned* matrices. Section 14.8 develops the topic.

14.3 The eigenvalue itself

We can finally now approach the last major agent of matrix arithmetic, the *eigenvalue*. Suppose a square, $n \times n$ matrix A , a nonzero n -element vector

$$\mathbf{v} = I_n \mathbf{v} \neq 0, \quad (14.17)$$

and a scalar λ , together such that

$$A\mathbf{v} = \lambda\mathbf{v}, \quad (14.18)$$

or in other words such that $A\mathbf{v} = \lambda I_n \mathbf{v}$. If so, then

$$[A - \lambda I_n]\mathbf{v} = 0. \quad (14.19)$$

Since $I_n \mathbf{v}$ is nonzero, the last equation is true if and only if the matrix $[A - \lambda I_n]$ is singular—which in light of § 14.1.3 is to demand that

$$\det(A - \lambda I_n) = 0. \quad (14.20)$$

The left side of (14.20) is an n th-order polynomial in λ , the *characteristic polynomial*, whose n roots are the *eigenvalues*¹⁰ of the matrix A .

What is an eigenvalue, really? An eigenvalue is a scalar a matrix resembles under certain conditions. When a matrix happens to operate on the right *eigenvector* \mathbf{v} , it is all the same whether one applies the entire matrix or just the eigenvalue to the vector. The matrix scales the eigenvector by the eigenvalue without otherwise altering the vector, changing the vector's

¹⁰An example:

$$\begin{aligned} A &= \begin{bmatrix} 2 & 0 \\ 3 & -1 \end{bmatrix}, \\ \det(A - \lambda I_n) &= \det \begin{bmatrix} 2 - \lambda & 0 \\ 3 & -1 - \lambda \end{bmatrix} \\ &= (2 - \lambda)(-1 - \lambda) - (0)(3) \\ &= \lambda^2 - \lambda - 2 = 0, \\ \lambda &= -1 \text{ or } 2. \end{aligned}$$

magnitude but not its direction. The eigenvalue alone takes the place of the whole, hulking matrix. This is what (14.18) means. Of course it works only when \mathbf{v} happens to be the right *eigenvector*, which § 14.4 discusses.

When $\lambda = 0$, (14.20) makes $\det A = 0$, which as we have said is the sign of a singular matrix. Zero eigenvalues and singular matrices always travel together. *Singular matrices each have at least one zero eigenvalue; nonsingular matrices never do.*

The eigenvalues of a matrix's inverse are the inverses of the matrix's eigenvalues. That is,

$$\lambda'_j \lambda_j = 1 \quad \text{for all } 1 \leq j \leq n \text{ if } A'A = I_n = AA'. \quad (14.21)$$

The reason behind (14.21) comes from answering the question: if $A\mathbf{v}_j$ scales \mathbf{v}_j by the factor λ_j , then what does $A'A\mathbf{v}_j = I\mathbf{v}_j$ do to \mathbf{v}_j ?

Naturally one must solve (14.20)'s n th-order polynomial to locate the actual eigenvalues. One solves it by the same techniques by which one solves any polynomial: the quadratic formula (2.2); the cubic and quartic methods of Ch. 10; the Newton-Raphson iteration (4.30). On the other hand, the determinant (14.20) can be impractical to expand for a large matrix; here iterative techniques help: see [chapter not yet written].¹¹

14.4 The eigenvector

It is an odd fact that (14.19) and (14.20) reveal the eigenvalues λ of a square matrix A while obscuring the associated *eigenvectors* \mathbf{v} . Once one has calculated an eigenvalue, though, one can feed it back to calculate the associated eigenvector. According to (14.19), the eigenvectors are the n -element vectors for which

$$[A - \lambda I_n]\mathbf{v} = 0,$$

which is to say that the eigenvectors are the vectors of the kernel space of the degenerate matrix $[A - \lambda I_n]$ —which one can calculate (among other ways) by the Gauss-Jordan kernel formula (13.7) or the Gram-Schmidt kernel formula (13.59).

An eigenvalue and its associated eigenvector, taken together, are sometimes called an *eigensolution*.

¹¹The inexpensive [13] also covers the topic competently and readably.

14.5 Eigensolution facts

Many useful or interesting mathematical facts concern the eigensolution, among them the following.

- *If the eigensolutions of A are $(\lambda_j, \mathbf{v}_j)$, then the eigensolutions of $A + \alpha I_n$ are $(\lambda_j + \alpha, \mathbf{v}_j)$. The eigenvalues move over by αI_n while the eigenvectors remain fixed. This is seen by adding $\alpha \mathbf{v}_j$ to both sides of the definition $A\mathbf{v}_j = \lambda_j \mathbf{v}_j$.*
- *A matrix and its inverse share the same eigenvectors with inverted eigenvalues. Refer to (14.21) and its explanation in § 14.3.*
- *Eigenvectors corresponding to distinct eigenvalues are always linearly independent of one another. To prove this fact, consider several independent eigenvectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k-1}$ respectively with distinct eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_{k-1}$, and further consider another eigenvector \mathbf{v}_k which might or might not be independent but which too has a distinct eigenvalue λ_k . Were \mathbf{v}_k dependent, which is to say, did nontrivial coefficients c_j exist such that*

$$\mathbf{v}_k = \sum_{j=1}^{k-1} c_j \mathbf{v}_j,$$

then left-multiplying the equation by $A - \lambda_k$ would yield

$$0 = \sum_{j=1}^{k-1} (\lambda_j - \lambda_k) c_j \mathbf{v}_j,$$

impossible since the $k - 1$ eigenvectors are independent. Thus \mathbf{v}_k too is independent, whereupon by induction from a start case of $k = 1$ we conclude that there exists no dependent eigenvector with a distinct eigenvalue.

- *If an $n \times n$ square matrix A has n independent eigenvectors (which is always so if the matrix has n distinct eigenvalues and often so even otherwise), then any n -element vector can be expressed as a unique linear combination of the eigenvectors. This is a simple consequence of the fact that the $n \times n$ matrix V whose columns are the several eigenvectors \mathbf{v}_j has full rank $r = n$. Unfortunately some matrices with repeated eigenvalues also have repeated eigenvectors—as for example,*

curiously,¹² $\begin{bmatrix} 1 & 0; 1 & 1 \end{bmatrix}^T$, whose double eigenvalue $\lambda = 1$ has the single eigenvector $\begin{bmatrix} 1 & 0 \end{bmatrix}^T$. Section 14.10.2 speaks of matrices of the last kind.

- *An $n \times n$ square matrix whose eigenvectors are linearly independent of one another cannot share all eigensolutions with any other $n \times n$ square matrix.* This fact proceeds from the last point, that every n -element vector \mathbf{x} is a unique linear combination of independent eigenvectors. Neither of the two proposed matrices A_1 and A_2 could scale any of the eigenvector components of \mathbf{x} differently than the other matrix did, so $A_1\mathbf{x} - A_2\mathbf{x} = (A_1 - A_2)\mathbf{x} = 0$ for all \mathbf{x} , which in turn is possible only if $A_1 = A_2$.
- *A positive definite matrix has only real, positive eigenvalues. A nonnegative definite matrix has only real, nonnegative eigenvalues.* Were it not so, then $\mathbf{v}^*A\mathbf{v} = \lambda\mathbf{v}^*\mathbf{v}$ (in which $\mathbf{v}^*\mathbf{v}$ naturally is a positive real scalar) would violate the criterion for positive or nonnegative definiteness. See § 13.6.3.
- *Every $n \times n$ square matrix has at least one eigensolution if $n > 0$,* because according to the fundamental theorem of algebra (6.1) the matrix's characteristic polynomial (14.20) has at least one root, an eigenvalue, which by definition would be no eigenvalue if it had no eigenvector to scale, and for which (14.19) necessarily admits at least one nonzero solution \mathbf{v} because its matrix $A - \lambda I_n$ is degenerate.

14.6 Diagonalization

Any $n \times n$ matrix with n independent eigenvectors (which class per § 14.5 includes, but is not limited to, every $n \times n$ matrix with n distinct eigenvalues) can be *diagonalized* as

$$A = V\Lambda V^{-1}, \quad (14.22)$$

where

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 & 0 \\ 0 & \lambda_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & \lambda_{n-1} & 0 \\ 0 & 0 & \cdots & 0 & \lambda_n \end{bmatrix}$$

¹²[19]

is an otherwise empty $n \times n$ matrix with the eigenvalues of A set along its main diagonal and

$$V = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_{n-1} & \mathbf{v}_n \end{bmatrix}$$

is an $n \times n$ matrix whose columns are the eigenvectors of A . This is so because the identity $A\mathbf{v}_j = \mathbf{v}_j\lambda_j$ holds for all $1 \leq j \leq n$; or, expressed more concisely, because the identity

$$AV = V\Lambda$$

holds.¹³ The matrix V is invertible because its columns the eigenvectors are independent, from which (14.22) follows. Equation (14.22) is called the *eigenvalue decomposition*, the *diagonal decomposition* or the *diagonalization* of the square matrix A .

One might object that we had shown only how to compose some matrix $V\Lambda V^{-1}$ with the correct eigenvalues and independent eigenvectors, but had failed to show that the matrix was actually A . However, we need not show this, because § 14.5 has already demonstrated that two matrices with the same eigenvalues and independent eigenvectors are in fact the same matrix, whereby the product $V\Lambda V^{-1}$ can be nothing other than A .

An $n \times n$ matrix with n independent eigenvectors (which class, again, includes every $n \times n$ matrix with n distinct eigenvalues) is called a *diagonalizable* matrix. Besides factoring a diagonalizable matrix by (14.22), one can apply the same formula to compose a diagonalizable matrix with desired eigensolutions.

The diagonal matrix $\text{diag}\{\mathbf{x}\}$ of (11.55) is trivially diagonalizable as $\text{diag}\{\mathbf{x}\} = I_n \text{diag}\{\mathbf{x}\} I_n$.

It is a curious and useful fact that

$$A^2 = (V\Lambda V^{-1})(V\Lambda V^{-1}) = V\Lambda^2 V^{-1}$$

and by extension that

$$A^k = V\Lambda^k V^{-1} \tag{14.23}$$

for any diagonalizable matrix A . The diagonal matrix Λ^k is nothing more than the diagonal matrix Λ with each element individually raised to the k th power, such that

$$\left[\Lambda^k\right]_{ij} = \delta_{ij}\lambda_j^k.$$

¹³If this seems confusing, then consider that the j th column of the product AV is $A\mathbf{v}_j$, whereas the j th column of Λ having just the one element acts to scale V 's j th column only.

Changing $z \leftarrow k$ implies the generalization¹⁴

$$\begin{aligned} A^z &= V\Lambda^z V^{-1}, \\ \left[\Lambda^z\right]_{ij} &= \delta_{ij}\lambda_j^z, \end{aligned} \tag{14.24}$$

good for any diagonalizable A and complex z .

Nondiagonalizable matrices are troublesome and interesting. The nondiagonalizable matrix vaguely resembles the singular matrix in that both represent edge cases and can be hard to handle numerically; but the resemblance ends there, and a matrix can be either without being the other. The $n \times n$ null matrix for example is singular but still diagonalizable. What a nondiagonalizable matrix is in essence is a matrix with a repeated eigensolution: the same eigenvalue with the same eigenvector, twice or more. More formally, a nondiagonalizable matrix is a matrix with an n -fold eigenvalue whose corresponding eigenvector space fewer than n eigenvectors fully characterize. Section 14.10.2 will have more to say about the nondiagonalizable matrix.

14.7 Remarks on the eigenvalue

Eigenvalues and their associated eigenvectors stand among the principal reasons one goes to the considerable trouble to develop matrix theory as we have done in recent chapters. The idea that a matrix resembles a humble scalar in the right circumstance is powerful. Among the reasons for this is that a matrix can represent an iterative process, operating repeatedly on a vector \mathbf{v} to change it first to $A\mathbf{v}$, then to $A^2\mathbf{v}$, $A^3\mathbf{v}$ and so on. The *dominant eigenvalue* of A , largest in magnitude, tends then to transform \mathbf{v} into the associated eigenvector, gradually but relatively eliminating all other components of \mathbf{v} . Should the dominant eigenvalue have greater than unit magnitude, it destabilizes the iteration; thus one can sometimes judge the stability of a physical process indirectly by examining the eigenvalues of the matrix which describes it. Then there is the edge case of the nondiagonalizable matrix, which matrix surprisingly covers only part of its domain with eigenvectors. All this is fairly deep mathematics. It brings an appreciation of the matrix for reasons which were anything but apparent from the outset of Ch. 11.

Remarks continue in §§ 14.10.2 and 14.13.

¹⁴It may not be clear however according to (5.12) which branch of λ_j^z one should choose at each index j , especially if A has negative or complex eigenvalues.

14.8 Matrix condition

The largest in magnitude of the several eigenvalues of a diagonalizable operator A , denoted here λ_{\max} , tends to dominate the iteration $A^k \mathbf{x}$. Section 14.7 has named λ_{\max} the *dominant eigenvalue* for this reason.

One sometimes finds it convenient to normalize a dominant eigenvalue by defining a new operator $A' \equiv A/|\lambda_{\max}|$, whose own dominant eigenvalue $\lambda_{\max}/|\lambda_{\max}|$ has unit magnitude. In terms of the new operator, the iteration becomes $A^k \mathbf{x} = |\lambda_{\max}|^k A'^k \mathbf{x}$, leaving one free to carry the magnifying effect $|\lambda_{\max}|^k$ separately if one prefers to do so. However, the scale factor $1/|\lambda_{\max}|$ scales all eigenvalues equally; thus, if A 's eigenvalue of *smallest* magnitude is denoted λ_{\min} , then the corresponding eigenvalue of A' is $\lambda_{\min}/|\lambda_{\max}|$. If zero, then both matrices according to § 14.3 are singular; if nearly zero, then both matrices are ill conditioned.

Such considerations lead us to define the *condition* of a diagonalizable matrix quantitatively as¹⁵

$$\kappa \equiv \left| \frac{\lambda_{\max}}{\lambda_{\min}} \right|, \quad (14.25)$$

by which

$$\kappa \geq 1 \quad (14.26)$$

is always a real number of no less than unit magnitude. For best invertibility, $\kappa = 1$ would be ideal (it would mean that all eigenvalues had the same magnitude), though in practice quite a broad range of κ is usually acceptable. Could we always work in exact arithmetic, the value of κ might not interest us much as long as it stayed finite; but in computer floating point, or where the elements of A are known only within some tolerance, infinite κ tends to emerge imprecisely rather as large $\kappa \gg 1$. An *ill conditioned* matrix by definition¹⁶ is a matrix of large $\kappa \gg 1$. The applied mathematician handles such a matrix with due skepticism.

Matrix condition so defined turns out to have another useful application. Suppose that a diagonalizable matrix A is precisely known but that the corresponding driving vector \mathbf{b} is not. If

$$A(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} + \delta \mathbf{b},$$

¹⁵[7]

¹⁶There is of course no definite boundary, no particular edge value of κ , less than which a matrix is well conditioned, at and beyond which it turns ill conditioned; but you knew that already. If I tried to claim that a matrix with a fine $\kappa = 3$ were ill conditioned, for instance, or that one with a wretched $\kappa = 2^{0 \times 18}$ were well conditioned, then you might not credit me—but the mathematics nevertheless can only give the number; it remains to the mathematician to interpret it.

where $\delta \mathbf{b}$ is the error in \mathbf{b} and $\delta \mathbf{x}$ is the resultant error in \mathbf{x} , then one should like to bound the ratio $|\delta \mathbf{x}| / |\mathbf{x}|$ to ascertain the reliability of the solution \mathbf{x} . Transferring A to the equation's right side,

$$\mathbf{x} + \delta \mathbf{x} = A^{-1}(\mathbf{b} + \delta \mathbf{b}).$$

Subtracting $\mathbf{x} = A^{-1}\mathbf{b}$ and taking the magnitude,

$$|\delta \mathbf{x}| = |A^{-1} \delta \mathbf{b}|.$$

Dividing by $|\mathbf{x}| = |A^{-1}\mathbf{b}|$,

$$\frac{|\delta \mathbf{x}|}{|\mathbf{x}|} = \frac{|A^{-1} \delta \mathbf{b}|}{|A^{-1}\mathbf{b}|}.$$

The quantity $|A^{-1} \delta \mathbf{b}|$ cannot exceed $|\lambda_{\min}^{-1} \delta \mathbf{b}|$. The quantity $|A^{-1}\mathbf{b}|$ cannot fall short of $|\lambda_{\max}^{-1}\mathbf{b}|$. Thus,

$$\frac{|\delta \mathbf{x}|}{|\mathbf{x}|} \leq \frac{|\lambda_{\min}^{-1} \delta \mathbf{b}|}{|\lambda_{\max}^{-1}\mathbf{b}|} = \left| \frac{\lambda_{\max}}{\lambda_{\min}} \right| \frac{|\delta \mathbf{b}|}{|\mathbf{b}|}.$$

That is,

$$\frac{|\delta \mathbf{x}|}{|\mathbf{x}|} \leq \kappa \frac{|\delta \mathbf{b}|}{|\mathbf{b}|}. \quad (14.27)$$

Condition, incidentally, might technically be said to apply to scalars as well as to matrices, but ill condition remains a property of matrices alone. According to (14.25), the condition of every nonzero scalar is happily $\kappa = 1$.

14.9 The similarity transformation

Any collection of vectors assembled into a matrix can serve as a *basis* by which other vectors can be expressed. For example, if the columns of

$$B = \begin{bmatrix} 1 & -1 \\ 0 & 2 \\ 0 & 1 \end{bmatrix}$$

are regarded as a basis, then the vector

$$B \begin{bmatrix} 5 \\ 1 \end{bmatrix} = 5 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + 1 \begin{bmatrix} -1 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 2 \\ 1 \end{bmatrix}$$

is $(5, 1)$ in the basis B : five times the first basis vector plus once the second. The basis provides the units from which other vectors can be built.

Particularly interesting is the $n \times n$, invertible *complete basis* B , in which the n basis vectors are independent and address the same full space the columns of I_n address. If

$$\mathbf{x} = B\mathbf{u}$$

then \mathbf{u} represents \mathbf{x} in the basis B . Left-multiplication by B evidently converts out of the basis. Left-multiplication by B^{-1} ,

$$\mathbf{u} = B^{-1}\mathbf{x},$$

then does the reverse, converting into the basis. One can therefore convert any operator A to work within a complete basis B by the successive steps

$$\begin{aligned} A\mathbf{x} &= \mathbf{b}, \\ AB\mathbf{u} &= \mathbf{b}, \\ [B^{-1}AB]\mathbf{u} &= B^{-1}\mathbf{b}, \end{aligned}$$

by which the operator $B^{-1}AB$ is seen to be the operator A , only transformed to work within the basis^{17,18} B .

The conversion from A into $B^{-1}AB$ is called a *similarity transformation*. If B happens to be unitary (§ 13.11), then the conversion is also called a *unitary transformation*. The matrix $B^{-1}AB$ the transformation produces is said to be *similar* (or, if B is unitary, *unitarily similar*) to the matrix A . We have already met the similarity transformation in §§ 11.5 and 12.2. Now we have the theory to appreciate it properly.

Probably the most important property of the similarity transformation is that it alters no eigenvalues. That is, if

$$A\mathbf{x} = \lambda\mathbf{x},$$

then, by successive steps,

$$\begin{aligned} B^{-1}A(BB^{-1})\mathbf{x} &= \lambda B^{-1}\mathbf{x}, \\ [B^{-1}AB]\mathbf{u} &= \lambda\mathbf{u}. \end{aligned} \tag{14.28}$$

¹⁷The reader may need to ponder the basis concept a while to grasp it, but the concept is simple once grasped and little purpose would be served by dwelling on it here. Basically, the idea is that one can build the same vector from alternate building blocks, not only from the standard building blocks $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$, etc.—except that the right word for the relevant “building block” is *basis vector*. The books [20] and [27] introduce the basis more gently; one might consult one of those if needed.

¹⁸The professional matrix literature sometimes distinguishes by typeface between the matrix B and the basis \mathbf{B} its columns represent. Such semantical distinctions seem a little too fine for applied use, though. This book just uses B .

The eigenvalues of A and the similar $B^{-1}AB$ are the same for any square, $n \times n$ matrix A and any invertible, square, $n \times n$ matrix B .

14.10 The Schur decomposition

The *Schur decomposition* of an arbitrary, $n \times n$ square matrix A is

$$A = QU_SQ^*, \quad (14.29)$$

where Q is an $n \times n$ unitary matrix whose inverse, as for any unitary matrix (§ 13.11), is $Q^{-1} = Q^*$; and where U_S is a general upper triangular matrix which can have any values (even zeros) along its main diagonal. The Schur decomposition is slightly obscure, is somewhat tedious to derive and is of limited use in itself, but serves a theoretical purpose.¹⁹ We derive it here for this reason.

14.10.1 Derivation

Suppose that²⁰ (for some reason, which will shortly grow clear) we have a matrix B of the form

$$B = \begin{bmatrix} \ddots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \cdots & * & * & * & * & * & * & * & * & \cdots \\ \cdots & 0 & * & * & * & * & * & * & * & \cdots \\ \cdots & 0 & 0 & * & * & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & * & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 0 & * & * & * & * & \cdots \\ \cdots & 0 & 0 & 0 & 0 & * & * & * & * & \cdots \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}, \quad (14.30)$$

¹⁹The alternative is to develop the interesting but difficult *Jordan canonical form*, which for brevity's sake this chapter prefers to omit.

²⁰This subsection assigns various capital Roman letters to represent the several matrices and submatrices it manipulates. Its choice of letters except in (14.29) is not standard and carries no meaning elsewhere. The writer had to choose some letters and these are ones he chose.

This footnote mentions the fact because good mathematical style avoid assigning letters that already bear a conventional meaning in a related context (for example, this book avoids writing $A\mathbf{x} = \mathbf{b}$ as $T\mathbf{e} = \mathbf{i}$, not because the latter is wrong but because it would be extremely confusing). The Roman alphabet provides only twenty-six capitals, though, of which this subsection uses too many to be allowed to reserve any. See Appendix B.

which contains a smaller, $(n-i) \times (n-i)$ unitary submatrix W_o in its lower right corner and resembles I_n elsewhere. Beginning from (14.31), we have by successive, reversible steps that

$$\begin{aligned} C &= W^*BW \\ &= (I_i + H_i W_o^* H_{-i})(B)(I_i + H_i W_o H_{-i}) \\ &= I_i B I_i + I_i B H_i W_o H_{-i} + H_i W_o^* H_{-i} B I_i \\ &\quad + H_i W_o^* H_{-i} B H_i W_o H_{-i}. \end{aligned}$$

The unitary submatrix W_o has only $n-i$ columns and $n-i$ rows, so $I_{n-i} W_o = W_o = W_o I_{n-i}$. Thus,

$$\begin{aligned} C &= I_i B I_i + I_i B H_i W_o I_{n-i} H_{-i} + H_i I_{n-i} W_o^* H_{-i} B I_i \\ &\quad + H_i I_{n-i} W_o^* I_{n-i} H_{-i} B H_i I_{n-i} W_o I_{n-i} H_{-i} \\ &= I_i [B] I_i + I_i [B H_i W_o H_{-i}] (I_n - I_i) + (I_n - I_i) [H_i W_o^* H_{-i} B] I_i \\ &\quad + (I_n - I_i) [H_i W_o^* B_o W_o H_{-i}] (I_n - I_i), \end{aligned}$$

where the last step has used (14.32) and the identity (11.76). The four terms on the equation's right, each term with rows and columns neatly truncated, represent the four quarters of $C \equiv W^*BW$ —upper left, upper right, lower left and lower right, respectively. The lower left term is null because

$$\begin{aligned} (I_n - I_i) [H_i W_o^* H_{-i} B] I_i &= (I_n - I_i) [H_i W_o^* I_{n-i} H_{-i} B I_i] I_i \\ &= (I_n - I_i) [H_i W_o^* H_{-i}] [(I_n - I_i) B I_i] I_i \\ &= (I_n - I_i) [H_i W_o^* H_{-i}] [0] I_i = 0, \end{aligned}$$

leaving

$$\begin{aligned} C &= I_i [B] I_i + I_i [B H_i W_o H_{-i}] (I_n - I_i) \\ &\quad + (I_n - I_i) [H_i W_o^* B_o W_o H_{-i}] (I_n - I_i). \end{aligned}$$

But the upper left term makes the upper left areas of B and C the same, and the upper right term does not bother us because we have not restricted the content of C 's upper right area. Apparently any $(n-i) \times (n-i)$ unitary submatrix W_o whatsoever obeys (14.31) in the lower left, upper left and upper right.

That leaves the lower right. Left- and right-multiplying (14.31) by the truncator $(I_n - I_i)$ to focus solely on the lower right area, we have the reduced requirement that

$$(I_n - I_i) C (I_n - I_i) = (I_n - I_i) W^* B W (I_n - I_i). \quad (14.34)$$

Further left-multiplying by H_{-i} , right-multiplying by H_i , and applying the identity (11.76) yields

$$I_{n-i}H_{-i}CH_iI_{n-i} = I_{n-i}H_{-i}W^*BW H_iI_{n-i};$$

or, substituting from (14.32),

$$C_o = I_{n-i}H_{-i}W^*BW H_iI_{n-i}.$$

Expanding W per (14.33),

$$C_o = I_{n-i}H_{-i}(I_i + H_iW_o^*H_{-i})B(I_i + H_iW_oH_{-i})H_iI_{n-i};$$

or, since $I_{n-i}H_{-i}I_i = 0 = I_iH_iI_{n-i}$,

$$\begin{aligned} C_o &= I_{n-i}H_{-i}(H_iW_o^*H_{-i})B(H_iW_oH_{-i})H_iI_{n-i} \\ &= I_{n-i}W_o^*H_{-i}BH_iW_oI_{n-i} \\ &= W_o^*I_{n-i}H_{-i}BH_iI_{n-i}W_o. \end{aligned}$$

Per (14.32), this is

$$C_o = W_o^*B_oW_o. \quad (14.35)$$

The steps from (14.34) to (14.35) are reversible, so the latter is as good a way to state the reduced requirement as the former is. To achieve a unitary transformation of the form (14.31), therefore, it suffices to satisfy (14.35).

The increasingly well-stocked armory of matrix theory we now have to draw from makes satisfying (14.35) possible as follows. Observe per § 14.5 that every square matrix has at least one eigensolution. Let $(\lambda_o, \mathbf{v}_o)$ represent an eigensolution of B_o —*any* eigensolution of B_o —with \mathbf{v}_o normalized to unit magnitude. Form the broad, $(n-i) \times (n-i+1)$ matrix

$$F \equiv \begin{bmatrix} \mathbf{v}_o & \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 & \cdots & \mathbf{e}_{n-i} \end{bmatrix}.$$

Decompose F by the Gram-Schmidt technique of § 13.10.2, choosing $p = 1$ during the first instance of the algorithm's step 3 (though choosing any permissible p thereafter), to obtain

$$F = Q_F R_F.$$

Noting that the Gram-Schmidt algorithm orthogonalizes only rightward, observe that the first column of the $(n-i) \times (n-i)$ unitary matrix Q_F remains simply the first column of F , which is the unit eigenvector \mathbf{v}_o :

$$[Q_F]_{*1} = Q_F \mathbf{e}_1 = \mathbf{v}_o.$$

Transform B_o unitarily by Q_F to define the new matrix

$$G \equiv Q_F^* B_o Q_F,$$

then transfer factors to reach the equation

$$Q_F G Q_F^* = B_o.$$

Right-multiplying by $Q_F \mathbf{e}_1 = \mathbf{v}_o$ and noting that $B_o \mathbf{v}_o = \lambda_o \mathbf{v}_o$, observe that

$$Q_F G \mathbf{e}_1 = \lambda_o \mathbf{v}_o.$$

Left-multiplying by Q_F^* ,

$$G \mathbf{e}_1 = \lambda_o Q_F^* \mathbf{v}_o.$$

Noting that the Gram-Schmidt process has rendered orthogonal to \mathbf{v}_o all columns of Q_F but the first, which is \mathbf{v}_o , observe that

$$G \mathbf{e}_1 = \lambda_o Q_F^* \mathbf{v}_o = \lambda_o \mathbf{e}_1 = \begin{bmatrix} \lambda_o \\ 0 \\ 0 \\ \vdots \end{bmatrix},$$

which means that

$$G = \begin{bmatrix} \lambda_o & * & * & \cdots \\ 0 & * & * & \cdots \\ 0 & * & * & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix},$$

which fits the very form (14.32) the submatrix C_o is required to have. Conclude therefore that

$$\begin{aligned} W_o &= Q_F, \\ C_o &= G, \end{aligned} \tag{14.36}$$

where Q_F and G as this paragraph develops, together constitute a valid choice for W_o and C_o , satisfying the reduced requirement (14.35) and thus also the original requirement (14.31).

Equation (14.36) completes a failsafe technique to transform unitarily any square matrix B of the form (14.30) into a square matrix C of the form (14.31). Naturally the technique can be applied recursively as

$$B|_{i=i'+1} = C|_{i=i'}, \quad 0 \leq i' \leq n-1, \tag{14.37}$$

because the form (14.30) of B at $i = i' + 1$ is nothing other than the form (14.31) of C at $i = i'$. Therefore, if we let

$$B|_{i=0} = A, \quad (14.38)$$

then it follows by induction that

$$B|_{i=n} = U_S, \quad (14.39)$$

where per (14.30) the matrix U_S has the general upper triangular form the Schur decomposition (14.29) requires. Moreover, because the product of unitary matrices according to (13.62) is itself a unitary matrix, we have that

$$Q = \prod_{i'=0}^{n-1} (W|_{i=i'}), \quad (14.40)$$

which along with (14.39) accomplishes the Schur decomposition.

14.10.2 The nondiagonalizable matrix

The characteristic equation (14.20) of the general upper triangular matrix U_S is

$$\det(U_S - \lambda I_n) = 0.$$

Unlike most determinants, this determinant brings only the one term

$$\det(U_S - \lambda I_n) = \prod_{i=1}^n (u_{Sii} - \lambda) = 0$$

whose factors run straight down the main diagonal, where the determinant's $n! - 1$ other terms are all zero because each of them includes at least one zero factor from below the main diagonal.²¹ Hence no element above the main diagonal of U_S even influences the eigenvalues, which apparently are

$$\lambda_i = u_{Sii}, \quad (14.41)$$

²¹The determinant's definition in § 14.1 makes the following two propositions equivalent: (i) that a determinant's term which includes one or more factors above the main diagonal also includes one or more factors below; (ii) that the only permutor that marks no position below the main diagonal is the one which also marks no position above. In either form, the proposition's truth might seem less than obvious until viewed from the proper angle.

Consider a permutor P . If P marked no position below the main diagonal, then it would necessarily have $p_{nn} = 1$, else the permutor's bottom row would be empty which is not allowed. In the next-to-bottom row, $p_{(n-1)(n-1)} = 1$, because the n th column is already occupied. In the next row up, $p_{(n-2)(n-2)} = 1$; and so on, thus affirming the proposition.

the main-diagonal elements.

Equation (14.28) preserves eigenvalues through a similarity transformation. The Schur decomposition (14.29) is in fact a similarity transformation; and, as we have seen, every matrix A has a Schur decomposition. If therefore

$$A = QU_SQ^*,$$

then *the eigenvalues of A are just the values along the main diagonal of U_S* .²²

One might think that the Schur decomposition offered an easy way to calculate eigenvalues; but it is less easy than it first appears because one must calculate eigenvalues to reach the Schur decomposition in the first place. Whatever practical merit the Schur decomposition might have or lack, however, it brings at least the theoretical benefit of (14.41): every square matrix without exception has a Schur decomposition, whose triangular factor U_S openly lists all eigenvalues along its main diagonal.

This theoretical benefit pays when some of the n eigenvalues of an $n \times n$ square matrix A repeat. By the Schur decomposition, one can construct a second square matrix A' , as near as desired to A but having n distinct eigenvalues, simply by perturbing the main diagonal of U_S to²³

$$\begin{aligned} U'_S &\equiv U_S + \epsilon \operatorname{diag}\{\mathbf{u}\}, \\ u_{i'} &\neq u_i \text{ if } \lambda_{i'} = \lambda_i, \end{aligned} \tag{14.42}$$

where $|\epsilon| \ll 1$ and where \mathbf{u} is an arbitrary vector that meets the criterion given. Though infinitesimally near A , the modified matrix $A' = QU'_SQ^*$ unlike A has n (maybe infinitesimally) distinct eigenvalues. With sufficient toil, one can analyze such perturbed eigenvalues and their associated eigenvectors similarly as § 9.6.2 has analyzed perturbed poles.

²²An unusually careful reader might worry that A and U_S had the same eigenvalues with different multiplicities. It would be surprising if it actually were so; but, still, one would like to give a sounder reason than the participle “surprising.” Consider however that

$$\begin{aligned} A - \lambda I_n &= QU_SQ^* - \lambda I_n = Q[U_S - Q^*(\lambda I_n)Q]Q^* \\ &= Q[U_S - \lambda(Q^*I_nQ)]Q^* = Q[U_S - \lambda I_n]Q^*. \end{aligned}$$

According to (14.11) and (14.13), this equation’s determinant is

$$\det[A - \lambda I_n] = \det\{Q[U_S - \lambda I_n]Q^*\} = \det Q \det[U_S - \lambda I_n] \det Q^* = \det[U_S - \lambda I_n],$$

which says that A and U_S have not only the same eigenvalues but also the same characteristic polynomials, thus further the same multiplicities of eigenvalues.

²³Equation (11.55) defines the $\operatorname{diag}\{\cdot\}$ notation.

Equation (14.42) brings us to the nondiagonalizable matrix of the subsection's title. Section 14.6 and its diagonalization formula (14.22) diagonalize any matrix with distinct eigenvalues and even any matrix with repeated eigenvalues but distinct eigenvectors, but fail where eigenvectors repeat. Equation (14.42) separates eigenvalues, thus also eigenvectors—for according to § 14.5 eigenvectors of distinct eigenvalues never depend on one another—permitting a nonunique but still sometimes usable form of diagonalization in the limit $\epsilon \rightarrow 0$ even when the matrix in question is nondiagonalizable.

The finding that every matrix is arbitrarily nearly diagonalizable illuminates a question the chapter has evaded up to the present point. The question: does a p -fold root in the characteristic polynomial (14.20) necessarily imply a p -fold eigenvalue in the corresponding matrix? The existence of the nondiagonalizable matrix casts a shadow of doubt until one realizes that every nondiagonalizable matrix is arbitrarily nearly diagonalizable—and, better, is arbitrarily nearly diagonalizable with distinct eigenvalues. If you claim that a matrix has a triple eigenvalue and someone disputes the claim, then you can show him a nearly identical matrix with three infinitesimally distinct eigenvalues. That is the essence of the idea. We shall leave the answer in that form.

Generalizing the nondiagonalizability concept leads one eventually to the ideas of the *generalized eigenvector*²⁴ (which solves the higher-order linear system $[A - \lambda I]^k \mathbf{v} = 0$) and the *Jordan canonical form*,²⁵ which together roughly track the sophisticated conventional pole-separation technique of § 9.6.5. Then there is a kind of sloppy Schur form called a Hessenberg form which allows content in U_S along one or more subdiagonals just beneath the main diagonal. One could profitably propose and prove any number of useful theorems concerning the nondiagonalizable matrix and its generalized eigenvectors, or concerning the eigenvalue problem²⁶ more broadly, in more and less rigorous ways, but for the time being we shall let the matter rest there.

14.11 The Hermitian matrix

An $m \times m$ square matrix A that is its own adjoint,

$$A^* = A, \tag{14.43}$$

²⁴[15, Ch. 7]

²⁵[13, Ch. 5]

²⁶[47]

is called a *Hermitian* or *self-adjoint* matrix. Properties of the Hermitian matrix include that

- its eigenvalues are real,
- its eigenvectors corresponding to distinct eigenvalues lie orthogonal to one another, and
- it is unitarily diagonalizable (§§ 13.11 and 14.6) such that

$$A = V\Lambda V^*. \quad (14.44)$$

That the eigenvectors lie orthogonal to one another is proved by letting (λ, \mathbf{v}) represent an eigensolution of A and constructing the product $\mathbf{v}^* A \mathbf{v}$, for which

$$\lambda^* \mathbf{v}^* \mathbf{v} = (A \mathbf{v})^* \mathbf{v} = \mathbf{v}^* A \mathbf{v} = \mathbf{v}^* (A \mathbf{v}) = \lambda \mathbf{v}^* \mathbf{v}.$$

That is,

$$\lambda^* = \lambda,$$

which naturally is possible only if λ is real.

That eigenvectors corresponding to distinct eigenvalues lie orthogonal to one another is proved²⁷ by letting $(\lambda_1, \mathbf{v}_1)$ and $(\lambda_2, \mathbf{v}_2)$ represent eigensolutions of A and constructing the product $\mathbf{v}_2^* A \mathbf{v}_1$, for which

$$\lambda_2^* \mathbf{v}_2^* \mathbf{v}_1 = (A \mathbf{v}_2)^* \mathbf{v}_1 = \mathbf{v}_2^* A \mathbf{v}_1 = \mathbf{v}_2^* (A \mathbf{v}_1) = \lambda_1 \mathbf{v}_2^* \mathbf{v}_1.$$

That is,

$$\lambda_2^* = \lambda_1 \quad \text{or} \quad \mathbf{v}_2^* \mathbf{v}_1 = 0.$$

But according to the last paragraph all eigenvalues are real; the eigenvalues λ_1 and λ_2 are no exceptions. Hence,

$$\lambda_2 = \lambda_1 \quad \text{or} \quad \mathbf{v}_2^* \mathbf{v}_1 = 0.$$

To prove the last hypothesis of the three needs first some definitions as follows. Given an $m \times m$ matrix A , let the s columns of the $m \times s$ matrix V_o represent the s independent eigenvectors of A such that (i) each column has unit magnitude and (ii) columns whose eigenvectors share the same eigenvalue lie orthogonal to one another. Let the $s \times s$ diagonal matrix Λ_o carry the eigenvalues on its main diagonal such that

$$A V_o = V_o \Lambda_o,$$

²⁷[27, § 8.1]

where the distinction between the matrix Λ_o and the full eigenvalue matrix Λ of (14.22) is that the latter always includes a p -fold eigenvalue p times, whereas the former includes a p -fold eigenvalue only as many times as the eigenvalue enjoys independent eigenvectors. Let the $m-s$ columns of the $m \times (m-s)$ matrix V_o^\perp represent the complete orthogonal complement (§ 13.9) to V_o —perpendicular to all eigenvectors, each column of unit magnitude—such that

$$V_o^{\perp*} V_o = 0 \quad \text{and} \quad V_o^{\perp*} V_o^\perp = I_{m-s}.$$

Recall from § 14.5 that $s \neq 0$ but $0 < s \leq m$ because every square matrix has at least one eigensolution. Recall from § 14.6 that $s = m$ if and only if A is diagonalizable.²⁸

With these definitions in hand, we can now prove that all Hermitian matrices are diagonalizable by contradiction, falsely supposing a nondiagonalizable Hermitian matrix A , whose V_o^\perp (since A is supposed to be nondiagonalizable, implying that $s < m$) would have at least one column. For such a matrix A , $s \times (m-s)$ and $(m-s) \times (m-s)$ auxiliary matrices F and G necessarily would exist such that

$$AV_o^\perp = V_o F + V_o^\perp G,$$

not due to any unusual property of the product AV_o^\perp but for the mundane reason that the columns of V_o and V_o^\perp together by definition addressed

²⁸A concrete example: the invertible but nondiagonalizable matrix

$$A = \begin{bmatrix} -1 & 0 & 0 & 0 \\ -6 & 5 & \frac{5}{2} & -\frac{5}{2} \\ 0 & 0 & 5 & 0 \\ 0 & 0 & 0 & 5 \end{bmatrix}$$

has a single eigenvalue at $\lambda = -1$ and a triple eigenvalue at $\lambda = 5$, the latter of whose eigenvector space is fully characterized by two eigenvectors rather than three such that

$$V_o = \begin{bmatrix} \frac{1}{\sqrt{2}} & 0 & 0 \\ \frac{1}{\sqrt{2}} & 1 & 0 \\ 0 & 0 & \frac{1}{\sqrt{2}} \\ 0 & 0 & \frac{1}{\sqrt{2}} \end{bmatrix}, \quad \Lambda_o = \begin{bmatrix} -1 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 5 \end{bmatrix}, \quad V_o^\perp = \begin{bmatrix} 0 \\ 0 \\ \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \end{bmatrix}.$$

The orthogonal complement V_o^\perp supplies the missing vector, not an eigenvector but perpendicular to them all.

In the example, $m = 4$ and $s = 3$.

All vectors in the example are reported with unit magnitude. The two $\lambda = 5$ eigenvectors are reported in mutually orthogonal form, but notice that eigenvectors corresponding to distinct eigenvalues need not be orthogonal when A is not Hermitian.

the space of *all* m -element vectors—including the columns of AV_o^\perp . Left-multiplying by V_o^* , we would have by successive steps that

$$\begin{aligned} V_o^* AV_o^\perp &= V_o^* V_o F + V_o^* V_o^\perp G, \\ (AV_o)^* V_o^\perp &= I_s F + V_o^* V_o^\perp G, \\ (V_o \Lambda_o)^* V_o^\perp &= F + V_o^* V_o^\perp G, \\ \Lambda_o^* V_o^* V_o^\perp &= F + V_o^* V_o^\perp G, \\ \Lambda_o^*(0) &= F + (0)G, \\ 0 &= F, \end{aligned}$$

where we had relied on the assumption that A were Hermitian and thus that, as proved above, its distinct eigenvectors lay orthogonal to one another; in consequence of which $A^* = A$ and $V_o^* V_o = I_s$.

The finding that $F = 0$ reduces the AV_o^\perp equation above to read

$$AV_o^\perp = V_o^\perp G.$$

In the reduced equation the matrix G would have at least one eigensolution, not due to any unusual property of G but because according to § 14.5 every square matrix, 1×1 or larger, has at least one eigensolution. Let (μ, \mathbf{w}) represent an eigensolution of G . Right-multiplying by the $(m - s)$ -element vector $\mathbf{w} \neq 0$, we would have by successive steps that

$$\begin{aligned} AV_o^\perp \mathbf{w} &= V_o^\perp G \mathbf{w}, \\ A(V_o^\perp \mathbf{w}) &= \mu(V_o^\perp \mathbf{w}). \end{aligned}$$

The last equation claims that $(\mu, V_o^\perp \mathbf{w})$ were an eigensolution of A , when we had supposed that all of A 's eigenvectors lay in the space addressed by the columns of V_o . The contradiction proves false the assumption that gave rise to it. The assumption: that a nondiagonalizable Hermitian A existed. We conclude that that all Hermitian matrices are diagonalizable—and conclude further that they are *unitarily* diagonalizable on the ground that their eigenvectors lie orthogonal to one another—as was to be demonstrated.

Having proven that all Hermitian matrices are diagonalizable and have real eigenvalues and orthogonal eigenvectors, one wonders whether the converse holds: are all diagonalizable matrices with real eigenvalues and orthogonal eigenvectors Hermitian? To show that they are, one can construct the matrix described by the diagonalization formula (14.22),

$$A = V \Lambda V^*,$$

where $V^{-1} = V^*$ because this V is unitary (§ 13.11). The equation's adjoint is

$$A^* = V\Lambda^*V^*.$$

But all the eigenvalues here are real, which means that $\Lambda^* = \Lambda$ and the right sides of the two equations are the same. As was to be demonstrated, therefore, $A^* = A$. *All diagonalizable matrices with real eigenvalues and orthogonal eigenvectors are Hermitian.*

This section brings properties that greatly simplify many kinds of matrix analysis. The properties demand a Hermitian matrix, which might seem a severe and unfortunate restriction—except that one can left-multiply any exactly determined linear system $C\mathbf{x} = \mathbf{d}$ by C^* to get the equivalent Hermitian system

$$[C^*C]\mathbf{x} = [C^*\mathbf{d}], \quad (14.45)$$

in which $A = C^*C$ and $\mathbf{b} = C^*\mathbf{d}$, for which the properties obtain.

14.12 The singular value decomposition

[This section is a new draft. The author requests error reports to thb@derivations.org if the date on the book's title page is less than a year old.]

Occasionally an elegant idea awaits discovery, overlooked, almost in plain sight. If the unlikely thought occurred to you to take the square root of a matrix, then the following idea is one you might discover.

Consider the $n \times n$ product A^*A of a tall or square, $m \times n$ matrix A of full column rank

$$r = n \leq m$$

and its adjoint A^* . The product A^*A is invertible according to § 13.6.2; is positive definite according to § 13.6.3; and, since $(A^*A)^* = A^*A$, is clearly Hermitian according to § 14.11; thus is unitarily diagonalizable according to (14.44) as

$$A^*A = V\Lambda V^*. \quad (14.46)$$

Here, the $n \times n$ matrices Λ and V represent the eigenvalues and eigenvectors not of A but of A^*A . Because A^*A is positive definite, all of its eigenvalues per § 14.5 are real and positive (which means among other things that Λ has

full rank). Letting the $n \times n$ *singular value* matrix $\Sigma = \sqrt{\Lambda}$ be such that

$$\begin{aligned} \Lambda &= \Sigma^* \Sigma, \\ \Sigma^* = \Sigma &= \begin{bmatrix} +\sqrt{\lambda_1} & 0 & \cdots & 0 & 0 \\ 0 & +\sqrt{\lambda_2} & \cdots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & +\sqrt{\lambda_{n-1}} & 0 \\ 0 & 0 & \cdots & 0 & +\sqrt{\lambda_n} \end{bmatrix} \end{aligned} \quad (14.47)$$

and applying (14.47) to (14.46) yields

$$\begin{aligned} A^* A &= V \Sigma^* \Sigma V^*, \\ V^* A^* A V &= \Sigma^* \Sigma. \end{aligned} \quad (14.48)$$

Now consider the $m \times m$ matrix U such that

$$\begin{aligned} AV \Sigma^{-1} &= U I_n, \\ AV &= U \Sigma, \\ A &= U \Sigma V^*. \end{aligned} \quad (14.49)$$

Substituting (14.49)'s second line into (14.48)'s second line gives the equation

$$\Sigma^* U^* U \Sigma = \Sigma^* \Sigma;$$

but $\Sigma \Sigma^{-1} = I_n$, so left- and right-multiplying respectively by Σ^{-*} and Σ^{-1} leaves that

$$I_n U^* U I_n = I_n,$$

which says neither more nor less than that the first n columns of U are orthonormal. Equation (14.49) does not constrain the last $m - n$ columns of U , leaving us free to make them anything we want. Why not use Gram-Schmidt to make them orthonormal, too, thus making U a unitary matrix? If we do this, then the surprisingly simple (14.49) constitutes the *singular value decomposition* of A .

If A happens to have broad shape then we can decompose A^* , instead, so this case poses no special trouble. Neither does the matrix of less than full column rank. Apparently every full-rank matrix has a singular value decomposition.

But what of the matrix of less than full rank? In this case the product $A^* A$ is not invertible and has only r nonzero eigenvalues. However, if the nonzero eigenvalues are arranged first in Λ , then (14.49) becomes

$$\begin{aligned} AV \Sigma^{-1} &= U I_r, \\ AV &= U \Sigma, \\ A &= U \Sigma V^*. \end{aligned} \quad (14.50)$$

The product A^*A is nonnegative definite in this case and $\Sigma\Sigma^{-1} = I_r$, but the reasoning is otherwise the same as before. Apparently every matrix of less than full rank has a singular value decomposition, too.

If A is square and has full rank, then the singular value decomposition inverts the matrix as

$$A^{-1} = V\Sigma^{-1}U^*. \quad (14.51)$$

14.13 General remarks on the matrix

Chapters 11 through 14 have derived the uncomfortably bulky but—incredibly—approximately minimal knot of theory one needs to grasp the matrix properly and to use it with moderate versatility. As far as the writer knows, no one has yet discovered a satisfactory way to untangle the knot. The choice to learn the basic theory of the matrix is almost an all-or-nothing choice; and how many scientists and engineers would rightly choose the “nothing” if the matrix did not serve so very many applications as it does? Since it does serve so very many, the “all” it must be.²⁹ Applied mathematics brings nothing else quite like it.

These several matrix chapters have not covered every topic they might. The topics they omit fall roughly into two classes. One is the class of more advanced and more specialized matrix theory, about which we shall have more to say in a moment. The other is the class of basic matrix theory these chapters do not happen to use. The essential agents of matrix analysis—multiplicative associativity, rank, inversion, pseudoinversion, the kernel, the orthogonal complement, orthonormalization, the eigenvalue, diagonalization and so on—are the same in practically all books on the subject, but the way the agents are developed differs. This book has chosen a way that needs some tools like truncators other books omit, but does not need other tools like projectors other books³⁰ include. What has given these chapters their hefty bulk is not so much the immediate development of the essential agents as the preparatory development of theoretical tools used to construct the essential agents, yet most of the tools are of limited interest in themselves; it is the agents that matter. Tools like the projector not used here tend to be omitted here or deferred to later chapters, not because they are altogether useless but because they are not used *here* and because the present chapters are already

²⁹Of course, one might avoid understanding and work by memorized rules. That is not always a bad plan, really; but if that were *your* plan then it seems spectacularly unlikely that you would be reading a footnote buried beneath the further regions of the hinterland of Chapter 14 in such a book as this.

³⁰Such as [20, § 3.VI.3], a lengthy but well-knit tutorial this writer recommends.

too long. The reader who understands the Moore-Penrose pseudoinverse and/or the Gram-Schmidt process reasonably well can after all pretty easily figure out how to construct a projector without explicit instructions thereto, should the need arise.³¹

Paradoxically and thankfully, more advanced and more specialized matrix theory though often harder tends to come in smaller, more manageable increments: the Cholesky decomposition, for instance; or the conjugate-gradient algorithm. The theory develops endlessly. From the present pause one could proceed directly to such topics. However, since this is the *first* proper pause these several matrix chapters have afforded, since the book is *Derivations of Applied Mathematics* rather than *Derivations of Applied Matrices*, maybe we ought to take advantage to change the subject.

³¹Well, since we have brought it up (though only as an example of tools these chapters have avoided bringing up), briefly: a projector is a matrix that flattens an arbitrary vector \mathbf{b} into its nearest shadow $\tilde{\mathbf{b}}$ within some restricted subspace. If the columns of A represent the subspace, then \mathbf{x} represents $\tilde{\mathbf{b}}$ in the subspace basis iff $A\mathbf{x} = \tilde{\mathbf{b}}$, which is to say that $A\mathbf{x} \approx \mathbf{b}$, whereupon $\mathbf{x} = A^\dagger \mathbf{b}$. That is, per (13.32),

$$\tilde{\mathbf{b}} = A\mathbf{x} = AA^\dagger \mathbf{b} = [BC][C^*(CC^*)^{-1}(B^*B)^{-1}B^*]\mathbf{b} = B(B^*B)^{-1}B^*\mathbf{b},$$

in which the matrix $B(B^*B)^{-1}B^*$ is the projector. Thence it is readily shown that the deviation $\mathbf{b} - \tilde{\mathbf{b}}$ lies orthogonal to the shadow $\tilde{\mathbf{b}}$. One can approach the projector in other ways, but there is one way at least.

Chapter 15

Vector algebra (to be written)

This chapter on vector algebra remains to be typeset, first of a tentatively planned series of three on the mathematics of three-dimensional vector fields. Since this chapter is not here yet, in place of the chapter let us outline the rest of the book. Topics planned include the following:

- Vector algebra (this chapter);
- The divergence theorem and Stokes' theorem;
- Vector calculus;
- Probability (including the Gaussian, Rayleigh and Poisson distributions and the basic concept of sample variance);
- The gamma function;
- The Fourier transform (including the Gaussian pulse and the Laplace transform);
- Feedback control;
- The spatial Fourier transform;
- Transformations to speed series convergence (such as the Poisson sum formula, Mosig's summation-by-parts technique and the Watson transformation);

- The scalar wave (or Helmholtz) equation¹ $(\nabla^2 + k^2)\psi(\mathbf{r}) = f(\mathbf{r})$ and Green's functions;
- Weyl and Sommerfeld forms (that is, spatial Fourier transforms in which only two or even only one of the three dimensions is transformed) and the parabolic wave equation;
- Bessel (especially Hankel) functions;
- Numerical integration;
- Iterative techniques (especially conjugate gradient).

The goal is to complete approximately the above list, plus appropriate additional topics as the writing brings them to light. Beyond the goal, a few more topics are conceivable:

- Legendre and spherical Bessel functions;
- Stochastics;
- The mathematics of energy conservation;
- Statistical mechanics (however, the treatment of ideal-gas particle speeds probably goes in the earlier Probability chapter);
- Rotational dynamics;
- Kepler's laws.

Yet further developments, if any, are hard to foresee.²

¹As you may know, this is called the “scalar wave equation” because $\psi(\mathbf{r})$ and $f(\mathbf{r})$ are scalar functions, even though the independent variable \mathbf{r} they are functions of is a vector. Vector wave equations also exist and are extremely interesting, for example in electromagnetics. However, those are typically solved by (cleverly) transforming them into superpositions of scalar wave equations. The scalar wave equation merits a chapter for this reason among others.

²Any plans—I should say, any wishes—beyond the topics listed are no better than daydreams. However, for my own notes if for no other reason, plausible topics include the following: Hamiltonian mechanics; electromagnetics; the statics of materials; the mechanics of materials; fluid mechanics; advanced special functions; thermodynamics and the mathematics of entropy; quantum mechanics; electric circuits; information theory; statistics. Life should be so long, eh? Well, we shall see. Like most authors perhaps, I write in my spare time, the supply of which is necessarily limited and unpredictable. (Family responsibilities and other duties take precedence. My wife says to me, “You have a lot of chapters to write. It will take you a long time.” She understates the problem.) The book targets the list ending in iterative techniques as its actual goal.

The book belongs to the open-source tradition, which means that you as reader have a stake in it if you wish. If you have read the book, or a substantial fraction of it, as far as it has yet gone, then you can help to improve it. Check <http://www.derivations.org/> for the latest revision, then write me at thb@derivations.org. I would most expressly solicit your feedback on typos, misprints, false or missing symbols and the like; such errors only mar the manuscript, so no such correction is too small. On a higher plane, if you have found any part of the book unnecessarily confusing, please tell how so. On no particular plane, if you just want to tell me what you have done with your copy of the book, what you have learned from it, or how you have cited it, why, what author does not appreciate such feedback? Send it in.

If you find a part of the book insufficiently rigorous, then that is another matter. I do not discourage such criticism and would be glad to hear it, but this book may not be well placed to meet it (the book might compromise by including a footnote which briefly suggests the outline of a more rigorous proof, but it tries not to distract the narrative by formalities which do not serve applications). If you want to detail Hölder spaces and Galois theory, or whatever, then my response is likely to be that there is already a surfeit of fine professional mathematics books in print; this just isn't that kind of book. On the other hand, the book does intend to derive every one of its results adequately from an applied perspective; if it fails to do so in your view, then maybe you and I should discuss the matter. Finding the right balance is not always easy.

At the time of this writing, readers are downloading the book at the rate of about three thousand copies per year directly through derivations.org. Some fraction of those, plus others who have installed the book as a Debian package or have acquired the book through secondary channels, actually have read it; now you stand among them.

Write as appropriate. More to come.

THB

Appendix A

Hexadecimal and other notational matters

The importance of conventional mathematical notation is hard to overstate. Such notation serves two distinct purposes: it conveys mathematical ideas from writer to reader; and it concisely summarizes complex ideas on paper to the writer himself. Without the notation, one would find it difficult even to think clearly about the math; to discuss it with others, nearly impossible.

The right notation is not always found at hand, of course. New mathematical ideas occasionally find no adequate preestablished notation, when it falls to the discoverer and his colleagues to establish new notation to meet the need. A more difficult problem arises when old notation exists but is inelegant in modern use.

Convention is a hard hill to climb, and rightly so. Nevertheless, slavish devotion to convention does not serve the literature well; for how else can notation improve over time, if writers will not incrementally improve it? Consider the notation of the algebraist Girolamo Cardano in his 1539 letter to Tartaglia:

[T]he cube of one-third of the coefficient of the unknown is greater in value than the square of one-half of the number. [31]

If Cardano lived today, surely he would express the same thought in the form

$$\left(\frac{a}{3}\right)^3 > \left(\frac{x}{2}\right)^2.$$

Good notation matters.

Although this book has no brief to overhaul applied mathematical notation generally, it does seek to aid the honorable cause of notational evolution

in a few specifics. For example, the book sometimes treats 2π implicitly as a single symbol, so that (for instance) the quarter revolution or right angle is expressed as $2\pi/4$ rather than as the less evocative $\pi/2$.

As a single symbol, of course, 2π remains a bit awkward. One wants to introduce some new symbol $\xi = 2\pi$ thereto. However, it is neither necessary nor practical nor desirable to leap straight to notational Utopia in one great bound. It suffices in print to improve the notation incrementally. If this book treats 2π sometimes as a single symbol—if such treatment meets the approval of slowly evolving convention—then further steps, the introduction of new symbols ξ and such, can safely be left incrementally to future writers.

A.1 Hexadecimal numerals

Treating 2π as a single symbol is a small step, unlikely to trouble readers much. A bolder step is to adopt from the computer science literature the important notational improvement of the hexadecimal numeral. No incremental step is possible here; either we leap the ditch or we remain on the wrong side. In this book, we choose to leap.

Traditional decimal notation is unobjectionable for measured quantities like 63.7 miles, \$1.32 million or 9.81 m/s^2 , but its iterative tenfold structure meets little or no aesthetic support in mathematical theory. Consider for instance the decimal numeral 127, whose number suggests a significant idea to the computer scientist, but whose decimal notation does nothing to convey the notion of the largest signed integer storable in a byte. Much better is the base-sixteen hexadecimal notation 0x7F, which clearly expresses the idea of $2^7 - 1$. To the reader who is not a computer scientist, the aesthetic advantage may not seem immediately clear from the one example, but consider the decimal number 2,147,483,647, which is the largest signed integer storable in a standard thirty-two bit word. In hexadecimal notation, this is 0x7FFF FFFF, or in other words $2^{0\text{x}1\text{F}} - 1$. The question is: which notation more clearly captures the idea?

To readers unfamiliar with the hexadecimal notation, to explain very briefly: hexadecimal represents numbers not in tens but rather in sixteens. The rightmost place in a hexadecimal numeral represents ones; the next place leftward, sixteens; the next place leftward, sixteens squared; the next, sixteens cubed, and so on. For instance, the hexadecimal numeral 0x1357 means “seven, plus five times sixteen, plus thrice sixteen times sixteen, plus once sixteen times sixteen times sixteen.” In hexadecimal, the sixteen symbols 0123456789ABCDEF respectively represent the numbers zero through

fifteen, with sixteen being written 0x10.

All this raises the sensible question: why sixteen?¹ The answer is that sixteen is 2^4 , so hexadecimal (base sixteen) is found to offer a convenient shorthand for binary (base two, the fundamental, smallest possible base). Each of the sixteen hexadecimal digits represents a unique sequence of exactly four bits (binary digits). Binary is inherently theoretically interesting, but direct binary notation is unwieldy (the hexadecimal number 0x1357 is binary 0001 0011 0101 0111), so hexadecimal is written in proxy.

The conventional hexadecimal notation is admittedly a bit bulky and unfortunately overloads the letters A through F, letters which when set in italics usually represent coefficients not digits. However, the real problem with the hexadecimal notation is not in the notation itself but rather in the unfamiliarity with it. The reason it is unfamiliar is that it is not often encountered outside the computer science literature, but it is not encountered because it is not used, and it is not used because it is not familiar, and so on in a cycle. It seems to this writer, on aesthetic grounds, that this particular cycle is worth breaking, so this book uses the hexadecimal for integers larger than 9. If you have never yet used the hexadecimal system, it is worth your while to learn it. For the sake of elegance, at the risk of challenging entrenched convention, this book employs hexadecimal throughout.

Observe that in some cases, such as where hexadecimal numbers are arrayed in matrices, this book may omit the cumbersome hexadecimal prefix “0x.” Specific numbers with physical dimensions attached appear seldom in this book, but where they do naturally decimal not hexadecimal is used: $v_{\text{sound}} = 331 \text{ m/s}$ rather than the silly-looking $v_{\text{sound}} = 0\text{x}14\text{B m/s}$.

Combining the hexadecimal and 2π ideas, we note here for interest’s sake that

$$2\pi \approx 0\text{x}6.487\text{F}.$$

A.2 Avoiding notational clutter

Good applied mathematical notation is not cluttered. Good notation does not necessarily include every possible limit, qualification, superscript and

¹An alternative advocated by some eighteenth-century writers was twelve. In base twelve, one quarter, one third and one half are respectively written 0.3, 0.4 and 0.6. Also, the hour angles (§ 3.6) come in neat increments of $(0.06)(2\pi)$ in base twelve, so there are some real advantages to that base. Hexadecimal, however, besides having momentum from the computer science literature, is preferred for its straightforward proxy of binary.

subscript. For example, the sum

$$S = \sum_{i=1}^M \sum_{j=1}^N a_{ij}^2$$

might be written less thoroughly but more readably as

$$S = \sum_{i,j} a_{ij}^2$$

if the meaning of the latter were clear from the context.

When to omit subscripts and such is naturally a matter of style and subjective judgment, but in practice such judgment is often not hard to render. The balance is between showing few enough symbols that the interesting parts of an equation are not obscured visually in a tangle and a haze of redundant little letters, strokes and squiggles, on the one hand; and on the other hand showing enough detail that the reader who opens the book directly to the page has a fair chance to understand what is written there without studying the whole book carefully up to that point. Where appropriate, this book often condenses notation and omits redundant symbols.

Appendix B

The Greek alphabet

Mathematical experience finds the Roman alphabet to lack sufficient symbols to write higher mathematics clearly. Although not completely solving the problem, the addition of the Greek alphabet helps. See Table B.1.

When first seen in mathematical writing, the Greek letters take on a wise, mysterious aura. Well, the aura is fine—the Greek letters are pretty—but don't let the Greek letters throw you. They're just letters. We use

Table B.1: The Roman and Greek alphabets.

ROMAN							
<i>Aa</i>	Aa	<i>Gg</i>	Gg	<i>Mm</i>	Mm	<i>Tt</i>	Tt
<i>Bb</i>	Bb	<i>Hh</i>	Hh	<i>Nn</i>	Nn	<i>Uu</i>	Uu
<i>Cc</i>	Cc	<i>Ii</i>	Ii	<i>Oo</i>	Oo	<i>Vv</i>	Vv
<i>Dd</i>	Dd	<i>Jj</i>	Jj	<i>Pp</i>	Pp	<i>Ww</i>	Ww
<i>Ee</i>	Ee	<i>Kk</i>	Kk	<i>Qq</i>	Qq	<i>Xx</i>	Xx
<i>Ff</i>	Ff	<i>Ll</i>	Ll	<i>Rr</i>	Rr	<i>Yy</i>	Yy
				<i>Ss</i>	Ss	<i>Zz</i>	Zz
GREEK							
$A\alpha$	alpha	$H\eta$	eta	$N\nu$	nu	$T\tau$	tau
$B\beta$	beta	$\Theta\theta$	theta	$\Xi\xi$	xi	$\Upsilon\upsilon$	upsilon
$\Gamma\gamma$	gamma	$\text{I}\iota$	iota	$\text{O}\omicron$	omicron	$\Phi\phi$	phi
$\Delta\delta$	delta	$\text{K}\kappa$	kappa	$\Pi\pi$	pi	$\text{X}\chi$	chi
$\text{E}\epsilon$	epsilon	$\Lambda\lambda$	lambda	$\text{P}\rho$	rho	$\Psi\psi$	psi
$\text{Z}\zeta$	zeta	$\text{M}\mu$	mu	$\Sigma\sigma$	sigma	$\Omega\omega$	omega

them not because we want to be wise and mysterious¹ but rather because

¹Well, you can use them to be wise and mysterious if you want to. It's kind of fun, actually, when you're dealing with someone who doesn't understand math—if what you want is for him to go away and leave you alone. Otherwise, we tend to use Roman and Greek letters in various conventional ways: lower-case Greek for angles; capital Roman for matrices; e for the natural logarithmic base; f and g for unspecified functions; i , j , k , m , n , M and N for integers; P and Q for metasyntactic elements (the mathematical equivalents of `foo` and `bar`); t , T and τ for time; d , δ and Δ for change; A , B and C for unknown coefficients; J , Y and H for Bessel functions; etc. Even with the Greek, there are still not enough letters, so each letter serves multiple conventional roles: for example, i as an integer, an a-c electric current, or—most commonly—the imaginary unit, depending on the context. Cases even arise where a quantity falls back to an alternate traditional letter because its primary traditional letter is already in use: for example, the imaginary unit falls back from i to j where the former represents an a-c electric current.

This is not to say that any letter goes. If someone wrote

$$e^2 + \pi^2 = O^2$$

for some reason instead of the traditional

$$a^2 + b^2 = c^2$$

for the Pythagorean theorem, you would not find that person's version so easy to read, would you? Mathematically, maybe it doesn't matter, but the choice of letters is not a

we simply do not have enough Roman letters. An equation like

$$\alpha^2 + \beta^2 = \gamma^2$$

says no more than does an equation like

$$a^2 + b^2 = c^2,$$

after all. The letters are just different.

Incidentally, we usually avoid letters like the Greek capital H (eta), which looks just like the Roman capital H, even though H (eta) is an entirely proper member of the Greek alphabet. Mathematical symbols are useful to us only inasmuch as we can visually tell them apart.

matter of arbitrary utility only but also of convention, tradition and style: one of the early writers in a field has chosen some letter—who knows why?—then the rest of us follow. This is how it usually works.

When writing notes for your own personal use, of course, you can use whichever letter you want. Probably you will find yourself using the letter you are used to seeing in print, but a letter is a letter; any letter will serve. When unsure which letter to use, just pick one; you can always go back and change the letter in your notes later if you want.

Appendix C

Manuscript history

The book in its present form is based on various unpublished drafts and notes of mine, plus some of my wife Kristie's (née Hancock), going back to 1983 when I was fifteen years of age. What prompted the contest I can no longer remember, but the notes began one day when I challenged a high-school classmate to prove the quadratic formula. The classmate responded that he didn't need to prove the quadratic formula because the proof was in the class math textbook, then counterchallenged me to prove the Pythagorean theorem. Admittedly obnoxious (I was fifteen, after all) but not to be out-done, I whipped out a pencil and paper on the spot and started working. But I found that I could not prove the theorem that day.

The next day I did find a proof in the school library,¹ writing it down, adding to it the proof of the quadratic formula plus a rather inefficient proof of my own invention to the law of cosines. Soon thereafter the school's chemistry instructor happened to mention that the angle between the tetrahedrally arranged four carbon-hydrogen bonds in a methane molecule was 109° , so from a symmetry argument I proved that result to myself, too, adding it to my little collection of proofs. That is how it started.²

The book actually has earlier roots than these. In 1979, when I was twelve years old, my father bought our family's first eight-bit computer. The computer's built-in *BASIC* programming-language interpreter exposed

¹A better proof is found in § 2.10.

²Fellow gear-heads who lived through that era at about the same age might want to date me against the disappearance of the slide rule. Answer: in my country, or at least at my high school, I was three years too young to use a slide rule. The kids born in 1964 learned the slide rule; those born in 1965 did not. I wasn't born till 1967, so for better or for worse I always had a pocket calculator in high school. My family had an eight-bit computer at home, too, as we shall see.

functions for calculating sines and cosines of angles. The interpreter's manual included a diagram much like Fig. 3.1 showing what sines and cosines were, but it never explained how the computer went about calculating such quantities. This bothered me at the time. Many hours with a pencil I spent trying to figure it out, yet the computer's trigonometric functions remained mysterious to me. When later in high school I learned of the use of the Taylor series to calculate trigonometrics, into my growing collection of proofs the series went.

Five years after the Pythagorean incident I was serving the U.S. Army as an enlisted troop in the former West Germany. Although those were the last days of the Cold War, there was no shooting war at the time, so the duty was peacetime duty. My duty was in military signal intelligence, frequently in the middle of the German night when there often wasn't much to do. The platoon sergeant wisely condoned neither novels nor cards on duty, but he did let the troops read the newspaper after midnight when things were quiet enough. Sometimes I used the time to study my German—the platoon sergeant allowed this, too—but I owned a copy of Richard P. Feynman's *Lectures on Physics* [11] which I would sometimes read instead.

Late one night the battalion commander, a lieutenant colonel and West Point graduate, inspected my platoon's duty post by surprise. A lieutenant colonel was a highly uncommon apparition at that hour of a quiet night, so when that old man appeared suddenly with the sergeant major, the company commander and the first sergeant in tow—the last two just routed from their sleep, perhaps—surprise indeed it was. The colonel may possibly have caught some of my unlucky fellows playing cards that night—I am not sure—but me, he caught with my boots unpolished, reading the *Lectures*.

I snapped to attention. The colonel took a long look at my boots without saying anything, as stormclouds gathered on the first sergeant's brow at his left shoulder, then asked me what I had been reading.

"Feynman's *Lectures on Physics*, sir."

"Why?"

"I am going to attend the university when my three-year enlistment is up, sir."

"I see." Maybe the old man was thinking that I would do better as a scientist than as a soldier? Maybe he was remembering when he had had to read some of the *Lectures* himself at West Point. Or maybe it was just the singularity of the sight in the man's eyes, as though he were a medieval knight at bivouac who had caught one of the peasant levies, thought to be illiterate, reading Cicero in the original Latin. The truth of this, we shall never know. What the old man actually said was, "Good work, son. Keep

it up.”

The stormclouds dissipated from the first sergeant’s face. No one ever said anything to me about my boots (in fact as far as I remember, the first sergeant—who saw me seldom in any case—never spoke to me again). The platoon sergeant thereafter explicitly permitted me to read the *Lectures* on duty after midnight on nights when there was nothing else to do, so in the last several months of my military service I did read a number of them. It is fair to say that I also kept my boots better polished.

In Volume I, Chapter 6, of the *Lectures* there is a lovely introduction to probability theory. It discusses the classic problem of the “random walk” in some detail, then states without proof that the generalization of the random walk leads to the Gaussian distribution

$$p(x) = \frac{\exp(-x^2/2\sigma^2)}{\sigma\sqrt{2\pi}}.$$

For the derivation of this remarkable theorem, I scanned the book in vain. One had no Internet access in those days, but besides a well equipped gym the Army post also had a tiny library, and in one yellowed volume in the library—who knows how such a book got there?—I did find a derivation of the $1/\sigma\sqrt{2\pi}$ factor.³ The exponential factor, the volume did not derive. Several days later, I chanced to find myself in Munich with an hour or two to spare, which I spent in the university library seeking the missing part of the proof, but lack of time and unfamiliarity with such a German site defeated me. Back at the Army post, I had to sweat the proof out on my own over the ensuing weeks. Nevertheless, eventually I did obtain a proof which made sense to me. Writing the proof down carefully, I pulled the old high-school math notes out of my military footlocker (for some reason I had kept the notes and even brought them to Germany), dusted them off, and added to them the new Gaussian proof.

That is how it has gone. To the old notes, I have added new proofs from time to time; and although somehow I have misplaced the original high-school leaves I took to Germany with me, the notes have nevertheless grown with the passing years. After I had left the Army, married, taken my degree at the university, begun work as a building construction engineer, and started a family—when the latest proof to join my notes was a mathematical justification of the standard industrial construction technique for measuring the resistance-to-ground of a new building’s electrical grounding system—I was delighted to discover that Eric W. Weisstein had compiled

³The citation is now unfortunately long lost.

and published [45] a wide-ranging collection of mathematical results in a spirit not entirely dissimilar to that of my own notes. A significant difference remained, however, between Weisstein's work and my own. The difference was and is fourfold:

1. Number theory, mathematical recreations and odd mathematical names interest Weisstein much more than they interest me; my own tastes run toward math directly useful in known physical applications. The selection of topics in each body of work reflects this difference.
2. Weisstein often includes results without proof. This is fine, but for my own part I happen to like proofs.
3. Weisstein lists results encyclopedically, alphabetically by name. I organize results more traditionally by topic, leaving alphabetization to the book's index, that readers who wish to do so can coherently read the book from front to back.⁴
4. I have eventually developed an interest in the free-software movement, joining it as a Debian Developer [8]; and by these lights and by the standard of the Debian Free Software Guidelines (DFSG) [9], Weisstein's work is not free. No objection to non-free work as such is raised here, but the book you are reading *is* free in the DFSG sense.

A different mathematical reference, even better in some ways than Weisstein's and (if I understand correctly) indeed free in the DFSG sense, is emerging at the time of this writing in the on-line pages of the general-purpose encyclopedia Wikipedia [46]. Although Wikipedia remains generally uncitable⁵ and forms no coherent whole, it is laden with mathematical

⁴There is an ironic personal story in this. As children in the 1970s, my brother and I had a 1959 World Book encyclopedia in our bedroom, about twenty volumes. It was then a bit outdated (in fact the world had changed tremendously in the fifteen or twenty years following 1959, so the encyclopedia was more than a bit outdated) but the two of us still used it sometimes. Only years later did I learn that my father, who in 1959 was fourteen years old, had bought the encyclopedia with money he had earned delivering newspapers daily before dawn, *and then had read the entire encyclopedia, front to back*. My father played linebacker on the football team and worked a job after school, too, so where he found the time or the inclination to read an entire encyclopedia, I'll never know. Nonetheless, it does prove that even an encyclopedia can be read from front to back.

⁵Some "Wikipedians" do seem actively to be working on making Wikipedia authoritatively citable. The underlying philosophy and basic plan of Wikipedia admittedly tend to thwart their efforts, but their efforts nevertheless seem to continue to progress. We shall see. Wikipedia is a remarkable, monumental creation.

knowledge, including many proofs, which I have referred to more than a few times in the preparation of this text.

A book can follow convention or depart from it; yet, though occasional departure might render a book original, frequent departure seldom renders a book good. Whether this particular book is original or good, neither or both, is for the reader to tell, but in any case the book does both follow and depart. Convention is a peculiar thing: at its best, it evolves or accumulates only gradually, patiently storing up the long, hidden wisdom of generations past; yet herein arises the ancient dilemma. Convention, in all its richness, in all its profundity, can, sometimes, stagnate at a local maximum, a hillock whence higher ground is achievable not by gradual ascent but only by descent first—or by a leap. Descent risks a bog. A leap risks a fall. One ought not run such risks without cause, even in such an inherently unconservative discipline as mathematics.

Well, the book does risk. It risks one leap at least: it employs hexadecimal numerals.

This book is bound to lose at least a few readers for its unorthodox use of hexadecimal notation (“The first primes are 2, 3, 5, 7, 0xB, . . .”). Perhaps it will gain a few readers for the same reason; time will tell. I started keeping my own theoretical math notes in hex a long time ago; at first to prove to myself that I could do hexadecimal arithmetic routinely and accurately with a pencil, later from aesthetic conviction that it was the right thing to do. Like other applied mathematicians, I’ve several own private notations, and in general these are not permitted to burden the published text. The hex notation is not my own, though. It existed before I arrived on the scene, and since I know of no math book better positioned to risk its use, I have with hesitation and no little trepidation resolved to let this book use it. Some readers will approve; some will tolerate; undoubtedly some will do neither. The views of the last group must be respected, but in the meantime the book has a mission; and crass popularity can be only one consideration, to be balanced against other factors. The book might gain even more readers, after all, had it no formulas, and painted landscapes in place of geometric diagrams! I like landscapes, too, but anyway you can see where that line of logic leads.

More substantively: despite the book’s title, adverse criticism from some quarters for lack of rigor is probably inevitable; nor is such criticism necessarily improper from my point of view. Still, serious books by professional mathematicians tend to be *for* professional mathematicians, which is understandable but does not always help the scientist or engineer who wants to use the math to model something. The ideal author of such a book as this

would probably hold two doctorates: one in mathematics and the other in engineering or the like. The ideal author lacking, I have written the book.

So here you have my old high-school notes, extended over twenty years and through the course of two-and-a-half university degrees, now partly typed and revised for the first time as a \LaTeX manuscript. Where this manuscript will go in the future is hard to guess. Perhaps the revision you are reading is the last. Who can say? The manuscript met an uncommonly enthusiastic reception at Debconf 6 [8] May 2006 at Oaxtepec, Mexico; and in August of the same year it warmly welcomed Karl Sarnow and Xplora Knoppix [48] aboard as the second official distributor of the book. Such developments augur well for the book's future at least. But in the meantime, if anyone should challenge you to prove the Pythagorean theorem on the spot, why, whip this book out and turn to § 2.10. That should confound 'em.

THB

Bibliography

- [1] Larry C. Andrews. *Special Functions of Mathematics for Engineers*. McGraw-Hill, New York, 2nd edition, 1992.
- [2] Christopher Beattie, John Rossi, and Robert C. Rogers. *Notes on Matrix Theory*. Unpublished, Department of Mathematics, Virginia Polytechnic Institute and State University, Blacksburg, Va., 6 Dec. 2001.
- [3] Kristie H. Black. Private conversation, 1996.
- [4] Thaddeus H. Black. *Derivations of Applied Mathematics*. The Debian Project, <http://www.debian.org/>, 10 November 2007.
- [5] Gary S. Brown. Private conversation, 1996.
- [6] Richard Courant and David Hilbert. *Methods of Mathematical Physics*. Interscience (Wiley), New York, first English edition, 1953.
- [7] Eric de Sturler. Lecture, Virginia Polytechnic Institute and State University, Blacksburg, Va., 2007.
- [8] The Debian Project. <http://www.debian.org/>.
- [9] The Debian Project. Debian Free Software Guidelines, version 1.1. http://www.debian.org/social_contract#guidelines.
- [10] G. Doetsch. *Guide to the Applications of the Laplace and z-Transforms*. Van Nostrand Reinhold, London, 1971. Referenced indirectly by way of [32].
- [11] Richard P. Feynman, Robert B. Leighton, and Matthew Sands. *The Feynman Lectures on Physics*. Addison-Wesley, Reading, Mass., 1963–65. Three volumes.

- [12] Stephen D. Fisher. *Complex Variables*. Books on Mathematics. Dover, Mineola, N.Y., 2nd edition, 1990.
- [13] Joel N. Franklin. *Matrix Theory*. Books on Mathematics. Dover, Mineola, N.Y., 1968.
- [14] The Free Software Foundation. GNU General Public License, version 2. `/usr/share/common-licenses/GPL-2` on a Debian system. The Debian Project: <http://www.debian.org/>. The Free Software Foundation: 51 Franklin St., Fifth Floor, Boston, Mass. 02110-1301, USA.
- [15] Stephen H. Friedberg, Arnold J. Insel, and Lawrence E. Spence. *Linear Algebra*. Pearson Education/Prentice-Hall, Upper Saddle River, N.J., 4th edition, 2003.
- [16] Edward Gibbon. *The History of the Decline and Fall of the Roman Empire*. 1788.
- [17] William Goldman. *The Princess Bride*. Ballantine, New York, 1973.
- [18] Richard W. Hamming. *Methods of Mathematics Applied to Calculus, Probability, and Statistics*. Books on Mathematics. Dover, Mineola, N.Y., 1985.
- [19] Harvard University (author unknown). “Math 21B Review, Jan. 13, 2004”. http://www.math.harvard.edu/archive/21b_fall_03/final/21breview.pdf. As retrieved 8 Oct. 2007.
- [20] Jim Hefferon. *Linear Algebra*. Mathematics, St. Michael’s College, Colchester, Vt., 20 May 2006. (The book is free software, and besides is the best book of the kind the author of the book you hold has encountered. As of 3 Nov. 2007 at least, one can download it from <http://joshua.smcvt.edu/linearalgebra/>.)
- [21] Francis B. Hildebrand. *Advanced Calculus for Applications*. Prentice-Hall, Englewood Cliffs, N.J., 2nd edition, 1976.
- [22] Intel Corporation. *IA-32 Intel Architecture Software Developer’s Manual*, 19th edition, March 2006.
- [23] Brian W. Kernighan and Dennis M. Ritchie. *The C Programming Language*. Software Series. Prentice Hall PTR, Englewood Cliffs, N.J., 2nd edition, 1988.

- [24] Konrad Knopp. *Theory and Application of Infinite Series*. Hafner, New York, 2nd English ed., revised in accordance with the 4th German edition, 1947.
- [25] Werner E. Kohler. Lecture, Virginia Polytechnic Institute and State University, Blacksburg, Va., 2007.
- [26] The Labor Law Talk. http://encyclopedia.laborlawtalk.com/Applied_mathematics. As retrieved 1 Sept. 2005.
- [27] David C. Lay. *Linear Algebra and Its Applications*. Addison-Wesley, Reading, Mass., 1994.
- [28] N.N. Lebedev. *Special Functions and Their Applications*. Books on Mathematics. Dover, Mineola, N.Y., revised English edition, 1965.
- [29] David McMahon. *Quantum Mechanics Demystified*. Demystified Series. McGraw-Hill, New York, 2006.
- [30] Ali H. Nayfeh and Balakumar Balachandran. *Applied Nonlinear Dynamics: Analytical, Computational and Experimental Methods*. Series in Nonlinear Science. Wiley, New York, 1995.
- [31] John J. O'Connor and Edmund F. Robertson. *The MacTutor History of Mathematics*. School of Mathematics and Statistics, University of St. Andrews, Scotland, <http://www-history.mcs.st-andrews.ac.uk/>. As retrieved 12 Oct. 2005 through 2 Nov. 2005.
- [32] Charles L. Phillips and John M. Parr. *Signals, Systems and Transforms*. Prentice-Hall, Englewood Cliffs, N.J., 1995.
- [33] PlanetMath.org. *Planet Math*. <http://www.planetmath.org/>. As retrieved 21 Sept. 2007 through 20 Oct. 2007.
- [34] Carl Sagan. *Cosmos*. Random House, New York, 1980.
- [35] Wayne A. Scales. Lecture, Virginia Polytechnic Institute and State University, Blacksburg, Va., 2004.
- [36] Adel S. Sedra and Kenneth C. Smith. *Microelectronic Circuits*. Series in Electrical Engineering. Oxford University Press, New York, 3rd edition, 1991.
- [37] Al Shenk. *Calculus and Analytic Geometry*. Scott, Foresman & Co., Glenview, Ill., 3rd edition, 1984.

- [38] William L. Shirer. *The Rise and Fall of the Third Reich*. Simon & Schuster, New York, 1960.
- [39] Murray R. Spiegel. *Complex Variables: with an Introduction to Conformal Mapping and Its Applications*. Schaum's Outline Series. McGraw-Hill, New York, 1964.
- [40] Susan Stepney. "Euclid's proof that there are an infinite number of primes". <http://www-users.cs.york.ac.uk/susan/cyc/p/primeprf.htm>. As retrieved 28 April 2006.
- [41] James Stewart, Lothar Redlin, and Saleem Watson. *Precalculus: Mathematics for Calculus*. Brooks/Cole, Pacific Grove, Calif., 3rd edition, 1993.
- [42] Bjarne Stroustrup. *The C++ Programming Language*. Addison-Wesley, Boston, "special" (third-and-a-half?) edition, 2000.
- [43] Henk A. van der Vorst. *Iterative Krylov Methods for Large Linear Systems*. Number 13 in Monographs on Applied and Computational Mathematics. Cambridge University Press, Cambridge, 2003.
- [44] Eric W. Weisstein. *Mathworld*. <http://mathworld.wolfram.com/>. As retrieved 29 May 2006.
- [45] Eric W. Weisstein. *CRC Concise Encyclopedia of Mathematics*. Chapman & Hall/CRC, Boca Raton, Fla., 2nd edition, 2003.
- [46] The Wikimedia Foundation. *Wikipedia*. <http://en.wikipedia.org/>.
- [47] J.H. Wilkinson. *The Algebraic Eigenvalue Problem*. Monographs on Numerical Analysis. Clarendon Press, Oxford, 1965.
- [48] Xplora. *Xplora Knoppix*. <http://www.xplora.org/downloads/Knoppix/>.

Index

- ' , 51
- 0 (zero), 7
 - dividing by, 68
 - matrix, 244
 - vector, 244, 280
- 1 (one), 7, 45, 245
- 2π , 35, 45, 397
 - calculating, 185
- LU* decomposition, 282
- QR* decomposition, 348
 - inverting a matrix by, 355
- δ , 68, 146, 241
- ϵ , 68
- \equiv , 11
- \leftarrow , 10
- \ll and \gg , 69
- 0x, 398
- π , 397
- $d\ell$, 144
- dz , 168
- e , 91
- i , 39
- n*th root
 - calculation by Newton-Raphson, 88
- n*th-order expression, 219
- reductio ad absurdum*, 112, 307
- 16th century, 219
- absolute value, 40
- abstraction, 281
- accretion, 130
- active region, 244, 247
- addition
 - of rows or columns, 361
 - parallel, 118, 220
 - serial, 118
 - series, 118
- addition of matrices, 237
- addition operator
 - downward multitarget, 264
 - elementary, 251
 - leftward multitarget, 265
 - multitarget, 263
 - rightward multitarget, 265
 - upward multitarget, 264
- addition quasielementary, 263
- row, 264
- addressing a vector space, 302
- adjoint, 240, 361, 384
 - of a matrix inverse, 257
- aeronautical engineer, 331
- aileron, 331
- algebra
 - classical, 7
 - fundamental theorem of, 117
 - higher-order, 219
 - linear, 233
- algorithm
 - Gauss-Jordan, 288
 - implementation of from an equation, 348
- alternating signs, 175
- altitude, 38
- AMD, 290
- amortization, 198
- analytic continuation, 161
- analytic function, 43, 161
- angle, 35, 45, 55
 - double, 55
 - half, 55
 - hour, 55
 - interior, 35

- of a polygon, 36
 - of a triangle, 35
 - of rotation, 53
 - right, 45
 - square, 45
 - sum of, 35
- antiderivative, 130, 193
 - and the natural logarithm, 194
 - guessing, 197
 - of a product of exponentials, powers and logarithms, 216
- antiquity, 219
- applied mathematics, 1, 146, 156
 - foundations of, 233
- arc, 45
- arccosine, 45
 - derivative of, 108
- arcsine, 45
 - derivative of, 108
- arctangent, 45
 - derivative of, 108
- area, 7, 36, 139
 - surface, 139
- arg, 40
- Argand domain and range planes, 162
- Argand plane, 40
- Argand, Jean-Robert (1768–1822), 40
- arithmetic, 7
 - exact, 290, 305, 315, 334, 366, 374
 - of matrices, 237
- arithmetic mean, 122
- arithmetic series, 15
- arm, 98
- articles “a” and “the”, 333
- artillerist, 330
- assignment, 10
- associativity, 7
 - of matrix multiplication, 237
- average, 121
- axes, 51
 - changing, 51
 - rotation of, 51
- axiom, 2
- axis, 62
- azimuth, 330
- baroquity, 229
- basis, 375
 - complete, 376
 - converting to and from, 376
- battle, 330
- battlefield, 330
- binomial theorem, 72
- bit, 243, 290
- blackbody radiation, 29
- block, wooden, 70
- bond, 80
- borrower, 198
- bound on a power series, 175
- boundary condition, 198
- box, 7
- branch point, 39, 163
 - strategy to avoid, 164
- building construction, 335
- businessperson, 121
- C and C++, 9, 10
- calculus, 67, 125
 - fundamental theorem of, 130
 - the two complementary questions of, 67, 125, 131
- cannon, 330
- Cardano, Girolamo (also known as Cardanus or Cardan, 1501–1576), 219
- carriage wheel, 331
- Cauchy’s integral formula, 167, 200
- Cauchy, Augustin Louis (1789–1857), 167, 200
- chain rule, derivative, 81
- change, 67
- change of variable, 10
- change, rate of, 67
- characteristic polynomial, 368, 382, 384
- checking an integration, 143
- checking division, 143
- choice of wooden blocks, 70
- Cholesky, André-Louis (1875–1918), 391
- circle, 45, 55, 98
 - area of, 139
 - unit, 45

- circuit, 331
- cis, 102
- citation, 5
- city street, 332
- classical algebra, 7
- cleverness, 149, 227
- clock, 55
- closed analytic form, 216
- closed contour integration, 144
- closed form, 216
- closed surface integration, 142
- clutter, notational, 399
- coefficient, 233
 - inscrutable, 222
 - matching of, 26
 - nontrivial, 280, 370
 - unknown, 197
- coincident properties of the square matrix, 366
- column, 234, 291, 327, 350
 - addition of, 361
 - null, 360
 - null, appending a, 323
 - orthonormal, 354
 - scaled and repeated, 360
 - spare, 295, 315
- column operator, 239
- column rank, 338
 - full, 313, 315
- column vector, 307, 344
- combination, 70
 - properties of, 71
- combinatorics, 70
- commutation
 - hierarchy of, 253
 - of elementary operators, 253
 - of matrices, 254
 - of the identity matrix, 249
- commutivity, 7
 - noncommutivity of matrix multiplication, 237
- complete basis, 376
- completing the square, 12
- complex conjugation, 41
- complex contour
 - about a multiple pole, 173
- complex exponent, 100
- complex exponential, 91
 - and de Moivre's theorem, 101
 - derivative of, 103
 - inverse, derivative of, 103
 - properties, 104
- complex number, 5, 39, 65
 - actuality of, 108
 - conjugating, 41
 - imaginary part of, 41
 - is a scalar not a vector, 234
 - magnitude of, 40
 - multiplication and division, 41, 65
 - phase of, 40
 - real part of, 41
- complex plane, 40
- complex power, 74
- complex trigonometrics, 101
- complex variable, 5, 78
- composite number, 111
 - compositional uniqueness of, 112
- computer memory, 348
- computer processor, 290
- concert hall, 31
- condition, 367, 374
 - of a scalar, 375
- conditional convergence, 136
- cone
 - volume of, 139
- conjecture, 339
- conjugate, 39, 41, 109, 240
- conjugate transpose, 240
 - of a matrix inverse, 257
- conjugation, 41
- constant, 30
- constant expression, 11
- constant, indeterminate, 30
- constraint, 228
- contour, 145, 163
 - complex, 168, 173, 201
 - complex, about a multiple pole, 173
- contour integration, 144
 - closed, 144
 - closed complex, 167, 200

- complex, 168
 - of a vector quantity, 145
- contract, 121
- contractor, 335
- contradiction, proof by, 112, 307
- control, 331
- control surface, aeronautical, 331
- convention, 397
- convergence, 65, 135, 155
 - conditional, 136
 - domain of, 181
 - slow, 181
- coordinate rotation, 62
- coordinates, 62
 - cylindrical, 62
 - rectangular, 45, 62
 - relations among, 63
 - spherical, 62
- corner case, 226
- corner value, 220
- cosine, 45
 - derivative of, 103, 106
 - in complex exponential form, 102
 - law of cosines, 60
- counterexample, 237
- Courant, Richard (1888–1972), 3
- Cramer's rule, 366
- Cramer, Gabriel (1704–1752), 366
- cross-derivative, 191
- cross-term, 191
- cryptology, 111
- cubic expression, 11, 118, 219, 220
 - roots of, 223
- cubic formula, 223
- cubing, 224
- cylindrical coordinates, 62

- datum, 335
- day, 55
- de Moivre's theorem, 65
 - and the complex exponential, 101
- de Moivre, Abraham (1667–1754), 65
- Debian, 5
- Debian Free Software Guidelines, 5
- decomposition
 - LU , 282
 - QR , 348
 - diagonal, 371
 - differences between the Gram-Schmidt and Gauss-Jordan, 352
 - eigenvalue, 371
 - Gauss-Jordan, 282
 - Gram-Schmidt, 348
 - Gram-Schmidt, inverting a matrix by, 355
 - orthonormalizing, 348
 - Schur, 377
 - singular value, 388
- definite integral, 143
- definition, 2, 146
- definition notation, 11
- degenerate linear system, 313
- degenerate matrix, 312, 323
- degree of freedom, 330
- delta function, Dirac, 146
 - sifting property of, 146
- delta, Kronecker, 241
 - sifting property of, 241
- denominator, 22, 205
- density, 138
- dependent element, 327
- dependent variable, 76
- derivation, 1
- derivative, 67
 - balanced form, 75
 - chain rule for, 81
 - cross-, 191
 - definition of, 75
 - higher, 83
 - Jacobian, 276, 337, 343
 - Leibnitz notation for, 76
 - logarithmic, 80
 - manipulation of, 80
 - Newton notation for, 75
 - of z^a/a , 194
 - of a complex exponential, 103
 - of a function of a complex variable, 78
 - of a rational function, 210
 - of a trigonometric, 106

- of an inverse trigonometric, 108
 - of arcsine, arccosine and arctangent, 108
 - of sine and cosine, 103
 - of sine, cosine and tangent, 106
 - of the natural exponential, 92, 106
 - of the natural logarithm, 95, 108
 - of z^a , 79
 - product rule for, 81, 195, 277
 - second, 83
 - unbalanced form, 75
- determinant, 357
 - and the elementary operator, 362
 - definition of, 358
 - inversion by, 365
 - of a matrix inverse, 364
 - of a product, 364
 - of a unitary matrix, 364
 - product of two, 364
 - properties of, 359
 - rank- n , 358
 - zero, 363
- deviation, 340
- DFSG (Debian Free Software Guidelines), 5
- diag notation, the, 263
- diagonal, 36, 45
 - main, 242, 382
 - three-dimensional, 38
- diagonal decomposition, 371
- diagonal matrix, 262, 263
- diagonalizability, 372
- diagonalizable matrix, 385
- diagonalization, 371
- differentiability, 78
- differential equation, 198
 - solving by unknown coefficients, 198
- differentiation
 - analytical versus numeric, 144
- dimension, 234
- dimension-limited matrix, 244
- dimensionality, 51, 242, 316
- dimensionlessness, 45
- Dirac delta function, 146
 - sifting property of, 146
- Dirac, Paul (1902–1984), 146
- direction, 47
- discontinuity, 145
- distributivity, 7
- divergence to infinity, 38
- dividend, 22
- dividing by zero, 68
- division, 41, 65
 - by matching coefficients, 26
 - checking, 143
- divisor, 22
- domain, 38, 181, 327
 - sidestepping a, 182
- domain contour, 163
- domain neighborhood, 161
- dominant eigenvalue, 373, 374
- dot product, 344
- double angle, 55
- double integral, 138
- double pole, 206, 210
- double root, 225
- down, 45, 332
- downward multitarget addition operator, 264
- driving vector, 324, 327, 333
- dummy variable, 13, 133, 200
- east, 45, 330
- edge case, 5, 225, 280
- efficient implementation, 348
- eigensolution, 369, 370
 - count of, 371
 - impossibility to share given independent eigenvectors, 371
 - of a matrix inverse, 370
 - repeated, 373, 383
- eigenvalue, 357, 368
 - distinct, 370, 372
 - dominant, 373, 374
 - large, 374
 - magnitude of, 374
 - of a matrix inverse, 369
 - perturbed, 383
 - real, 385
 - repeated, 370, 383

- shifting, 370
 - small, 374
 - zero, 369, 374
- eigenvalue decomposition, 371
- eigenvector, 368, 369
 - generalized, 384
 - independent, 372
 - linearly independent, 370
 - of a matrix inverse, 370
 - orthogonal, 385
 - repeated, 373, 383
- electric circuit, 331
- electrical engineer, 331
- element, 234
 - free and dependent, 327
- elementary operator, 250
 - addition, 251, 256
 - and the determinant, 362
 - combining, 254
 - commutation of, 253
 - expanding, 254
 - interchange, 250, 255
 - inverse of, 250
 - invertibility of, 252
 - inverting, 254
 - scaling, 251, 256
 - sorting, 253
 - two, of different kinds, 253
- elementary similarity transformation, 282
- elementary vector, 249, 307, 344
 - left-over, 309
- elevation, 330
- elf, 39
- empty set, 280
- end, justifying the means, 306
- engineer, 331
- entire function, 165, 181, 187
- equation
 - simultaneous linear system of, 237, 324
 - solving a set of simultaneously, 53, 237, 324
 - superfluous, 324
- equator, 140
- error
 - due to rounding, 334
 - in the solution to a linear system, 374
- error bound, 175
- essential singularity, 39, 165, 190
- Euclid (c. 300 B.C.), 36, 112
- Euler's formula, 98
 - curious consequences of, 101
- Euler, Leonhard (1707–1783), 98, 100, 136
- evaluation, 83
- even function, 185
- evenness, 258
- exact arithmetic, 290, 305, 315, 334, 366, 374
- exact matrix, 305
- exact number, 305
- exact quantity, 305
- exactly determined linear system, 324
- exercises, 148
- existence, 212
- expansion of $1/(1 - z)^{n+1}$, 152
- expansion point
 - shifting, 157
- exponent, 16, 32
 - complex, 100
- exponent, floating-point, 290
- exponential, 32, 100
 - general, 94
 - integrating a product of a power and, 215
 - natural, error of, 180
 - resemblance of to x^∞ , 97
- exponential, complex, 91
 - and de Moivre's theorem, 101
- exponential, natural, 91
 - compared to x^a , 96
 - derivative of, 92, 106
 - existence of, 91
- exponential, real, 91
- extended operator, 243, 245, 355
 - decomposing, 301
- extension, 4
- extra column, 295, 315
- extremum, 83, 166

- factor
 - truncating, 298
- factorial, 13, 196
- factoring, 11
- factorization
 - QR , 348
 - full-rank, 314, 339, 352
 - Gauss-Jordan, 282
 - Gram-Schmidt, 348
 - orthonormalizing, 348
- factorization, prime, 112
- family of solutions, 332
- fast function, 96
- Ferrari, Lodovico (1522–1565), 219, 227
- flaw, logical, 114
- floating-point number, 290, 374
- forbidden point, 159, 162
- formalism, 134
- fourfold integral, 138
- Fourier transform
 - spatial, 138
- fraction, 205
- free element, 327
- freedom, degree of, 330
- freeway, 335
- frontier, 270
- Frullani's integral, 213
- Frullani, Giuliano (1795–1834), 213
- full column rank, 313, 315
- full rank, 312
- full rank factorization, 314
- full row rank, 313, 346
- full-rank factorization, 339, 352
- function, 38, 146
 - analytic, 43, 161
 - entire, 165, 181, 187
 - extremum of, 83
 - fast, 96
 - fitting of, 151
 - linear, 134
 - meromorphic, 165, 187
 - nonanalytic, 44
 - nonlinear, 134
 - odd or even, 185
 - of a complex variable, 78
 - rational, 205
 - rational, derivatives of, 210
 - rational, integral of, 209
 - single- and multiple-valued, 164
 - single- vs. multiple-valued, 162
 - slow, 96
- fundamental theorem of algebra, 117
- fundamental theorem of calculus, 130
- gamma function, 196
- Gauss, Carl Friedrich (1777–1855), 282, 326
- Gauss-Jordan algorithm, 288
- Gauss-Jordan decomposition, 282
 - and full column rank, 315
 - differences of against the Gram-Schmidt, 352
 - factors K and S of, 315
 - inverting the factors of, 297
- Gauss-Jordan factorization, 282
- Gauss-Jordan factors
 - properties of, 299
- Gauss-Jordan kernel formula, 326
- general exponential, 94
- general identity matrix, 245
- general interchange operator, 261
- General Public License, GNU, 5
- general scaling operator, 262
- general solution, 334
- general triangular matrix, 266, 382
- generality, 326
- generalized eigenvector, 384
- geometric majorization, 179
- geometric mean, 122
- geometric series, 29
 - majorization by, 179
- geometric series, variations on, 29
- geometrical arguments, 2
- geometrical visualization, 281
- geometry, 33
- GNU General Public License, 5
- GNU GPL, 5
- Goldman, William (1931–), 146
- GPL, 5
- Gram, Jørgen Pedersen (1850–1916), 346

- Gram-Schmidt decomposition, 348
 - differences of against the Gauss-Jordaniff, 119, 134, 155
 - 352
 - factor Q of, 352
 - factor S of, 352
 - inverting a matrix by, 355
- Gram-Schmidt kernel formula, 352
- Gram-Schmidt process, 346
- grapes, 109
- Greek alphabet, 401
- Greenwich, 55
- guessing roots, 229
- guessing the form of a solution, 197
- gunpowder, 330
- half angle, 55
- Hamming, Richard W. (1915–1998), 125, 156
- handwriting, reflected, 109
- harmonic mean, 122
- harmonic series, 179
- Heaviside unit step function, 145
- Heaviside, Oliver (1850–1925), 3, 109, 145
- Hermite, Charles (1822–1901), 240, 384
- Hermitian matrix, 384
- Hessenberg matrix, 384
- Hessenberg, Gerhard (1874–1925), 384
- hexadecimal, 398
- higher-order algebra, 219
- Hilbert, David (1862–1943), 3
- homogeneous solution, 333
- horse, 330
- hour, 55
- hour angle, 55
- hyperbolic functions, 102
 - properties of, 102
- hypotenuse, 36
- hypothesis, 328
- identity matrix, 245, 248
 - r -dimensional, 248
 - commutation of, 249
 - impossibility of promoting, 307
 - rank- r , 248
- identity, arithmetic, 7
- ill conditioned matrix, 367, 374
- imaginary number, 39
- imaginary part, 41
- imaginary unit, 39
- implementation
 - efficient, 348
- imprecise quantity, 305, 374
- indefinite integral, 143
- independence, 280, 312
- independent infinitesimal
 - variable, 130
- independent variable, 76
- indeterminate form, 85
- index, 236
- index of summation, 13
- induction, 41, 154
- inequality, 10, 340
- infinite differentiability, 161
- infinite dimensionality, 242
- infinitesimal, 68
 - and the Leibnitz notation, 76
 - dropping when negligible, 168
 - independent, variable, 130
 - practical size of, 68
 - second- and higher-order, 69
- infinity, 68
- inflection, 83
- inner product, 344
- integer, 15
 - composite, 111
 - compositional uniqueness of, 112
 - prime, 111
- integral, 125
 - and series summation, 178
 - as accretion or area, 126
 - as antiderivative, 130
 - as shortcut to a sum, 127
 - balanced form, 129
 - closed complex contour, 167, 200
 - closed contour, 144
 - closed surface, 142
 - complex contour, 168
 - concept of, 125

- contour, 144
- definite, 143
- double, 138
- fourfold, 138
- indefinite, 143
- multiple, 137
- of a rational function, 209
- sixfold, 138
- surface, 138, 140
- triple, 138
- vector contour, 145
- volume, 138
- integral swapping, 138
- integrand, 193
- integration
 - analytical versus numeric, 144
 - by antiderivative, 193
 - by closed contour, 200
 - by partial-fraction expansion, 205
 - by parts, 195
 - by substitution, 194
 - by Taylor series, 216
 - by unknown coefficients, 197
 - checking, 143
 - of a product of exponentials, powers and logarithms, 215
- integration techniques, 193
- Intel, 290
- interchange, 350, 359
 - refusing an, 315, 348
- interchange operator
 - elementary, 250
 - general, 261
- interchange quasidelementary, 261
- interest, 80, 198
- inverse
 - determinant of, 364
 - existence of, 322
 - mutual, 322
 - of a matrix product, 257
 - of a matrix transpose or adjoint, 257
 - rank- r , 320
 - uniqueness of, 322
- inverse complex exponential
 - derivative of, 103
- inverse trigonometric family of functions, 103
- inversion, 254, 319, 320
 - by determinant, 365
 - multiplicative, 7
 - symbolic, 357, 365
- inversion, arithmetic, 7
- invertibility, 338
 - of the elementary operator, 252
- irrational number, 115
- irreducibility, 2
- iteration, 86, 374
- Jacobi, Carl Gustav Jacob (1804–1851), 276
- Jacobian derivative, 276, 337, 343
- Jordan, Wilhelm (1842–1899), 282, 326
- kernel, 325
 - alternate formula for, 328
 - Gauss-Jordan formula for, 326
 - Gram-Schmidt formula for, 352
- kernel matrix, 325
 - converting between two of, 330
- kernel space, 326, 330
- Kronecker delta, 241
 - sifting property of, 241
- Kronecker, Leopold (1823–1891), 241
- l'Hôpital's rule, 84, 186
- l'Hôpital, Guillaume de (1661–1704), 84
- labor union, 335
- Laurent series, 21, 187
- Laurent, Pierre Alphonse (1813–1854), 21, 187
- law of cosines, 60
- law of sines, 59
- least squares, 335
- least-squares solution, 337
- leftward multitarget addition operator, 265
- leg, 36
- Leibnitz notation, 76, 130
- Leibnitz, Gottfried Wilhelm (1646–1716), 67, 76

- length
 - curved, 45
 - preservation of, 355
- limit, 69
- line, 51, 332
 - fitting a, 335
- linear algebra, 233
- linear combination, 134, 280, 312
- linear dependence, 280, 312
- linear expression, 11, 134, 219
- linear independence, 280, 312
- linear operator, 134
- linear superposition, 172
- linear system
 - classification of, 313
 - degenerate, 313
 - exactly determined, 324
 - nonoverdetermined, 332
 - nonoverdetermined, general solution of, 334
 - nonoverdetermined, particular solution of, 333
 - overdetermined, 313, 335, 336
 - taxonomy of, 313
 - underdetermined, 313
- linear transformation, 236
- linearity, 134
 - of a function, 134
 - of an operator, 134
- loan, 198
- locus, 117
- logarithm, 32
 - integrating a product of a power and, 215
 - natural, error of, 180
 - properties of, 33
 - resemblance of to x^0 , 97
- logarithm, natural, 94
 - and the antiderivative, 194
 - compared to x^a , 96
 - derivative of, 95, 108
- logarithmic derivative, 80
- logic, 306
 - reverse, 340
- logical flaw, 114
- lone-element matrix, 249
- long division, 22, 187
 - by $z - \alpha$, 116
 - procedure for, 25, 27
- loop, 349
- loop counter, 13
- lower triangular matrix, 265
- Maclaurin series, 160
- Maclaurin, Colin (1698–1746), 160
- magnification, 374
- magnitude, 40, 47, 65
 - of a vector, 345
 - of an eigenvalue, 374
 - preservation of, 355
 - unit, 374
- main diagonal, 242, 382
- majority, 349
- majorization, 156, 177, 178
 - geometric, 179
- maneuver, logical, 306
- mantissa, 290
- mapping, 38
- marking
 - permutor, 358
 - quasielementary, 358
- mason, 118
- mass density, 138
- matching coefficients, 26
- mathematician
 - applied, 1
 - professional, 2, 114, 212
- mathematics
 - applied, 1, 146, 156
 - professional or pure, 2, 114, 146
- mathematics, applied
 - foundations of, 233
- matrix, 233, 234
 - addition, 237
 - arithmetic of, 237
 - basic operations of, 236
 - basic use of, 236
 - broad, 312, 332
 - column of, 291, 327, 350
 - commutation of, 254

- degenerate, 312, 323
- diagonalizable, 385
- dimension-limited, 244
- dimensionality, 316
- exact, 305
- full-rank, 312
- general identity, 245
- Hermitian, 384
- identity, 245, 248
- identity, the impossibility of promoting, 307
- ill conditioned, 367, 374
- inversion of, 254, 319
- inversion properties of, 258
- inverting, 320, 365
- large, 334, 366
- lone-element, 249
- main diagonal of, 242, 382
- motivation for, 233
- motivation to, 236
- multiplication by a scalar, 237
- multiplication of, 237, 242
- multiplication, associativity of, 237
- multiplication, noncommutivity of, 237
- nondiagonalizable, 382
- nondiagonalizable versus singular, 373
- null, 244
- null column of, 350
- orthogonally complementary, 346
- padding with zeros, 242
- parallel unit triangular, properties of, 273
- perpendicular, 346
- projection, 390
- provenance of, 236
- raised to a complex power, 372
- rank of, 310
- rank- r inverse of, 258
- real, 337
- rectangular, 323
- row of, 291
- scalar, 245
- self-adjoint, 384
- singular, 323, 374
- sparse, 246
- square, 242, 312, 315, 320
- square, coincident properties of, 366
- tall, 312, 315
- triangular, construction of, 267
- truncating, 249
- unit lower triangular, 265
- unit parallel triangular, 269
- unit triangular, 265
- unit triangular, partial, 274
- unit upper triangular, 265
- unitary, 354
- matrix condition, 374
- matrix forms, 242
- matrix operator, 242
- matrix rudiments, 233, 279
- maximum, 83
- mean, 121
 - arithmetic, 122
 - geometric, 122
 - harmonic, 122
- means, justified by the end, 306
- memory, computer, 348
- meromorphic function, 165, 187
- minimum, 83
- minorization, 178
- mirror, 109
- model, 2, 114
- modulus, 40
- Moore, E.H. (1862–1932), 335
- Moore-Penrose pseudoinverse, 335, 352
- motivation, 324
- mountain road, 332
- multiple pole, 38, 206, 210
 - enclosing, 173
- multiple-valued function, 162, 164
- multiplication, 13, 41, 65
 - of matrices, 237, 242
- multiplicative inversion, 7
- multiplier, 233
- multitarget addition operator, 263
 - downward, 264
 - leftward, 265
 - rightward, 265
 - upward, 264

- multivariate Newton-Raphson iteration, 343
- Napoleon, 330
- natural exponential, 91
 - compared to x^a , 96
 - complex, 91
 - derivative of, 92, 106
 - error of, 180
 - existence of, 91
 - real, 91
- natural exponential family of functions, 103
- natural logarithm, 94
 - and the antiderivative, 194
 - compared to x^a , 96
 - derivative of, 95, 108
 - error of, 180
 - of a complex number, 101
- natural logarithmic family of functions, 103
- neighborhood, 161
- nesting, 349
- Newton, Sir Isaac (1642–1727), 67, 75, 86, 343
- Newton-Raphson iteration, 86, 219
 - multivariate, 343
- nobleman, 146
- nonanalytic function, 44
- nonanalytic point, 162, 171
- noncommutivity of matrix multiplication, 237
- nondiagonalizable matrix, 382
 - versus a singular matrix, 373
- noninvertibility, 322
- nonlinearity, 331
- nonnegative definiteness, 339, 371
- nonoverdetermined linear system, 332
- nonoverdetermined linear system, general solution of, 334
- nonoverdetermined linear system, particular solution of, 333
- normal vector or line, 139
- normalization, 346, 350, 374
- north, 45, 330
- null column, 350
- null matrix, 244
- null vector, 280
- number, 39
 - complex, 5, 39, 65
 - complex, actuality of, 108
 - exact, 305
 - imaginary, 39
 - irrational, 115
 - rational, 115
 - real, 39
 - very large or very small, 68
- number theory, 111
- numerator, 22, 205
- Observatory, Old Royal, 55
- Occam's razor
 - abusing, 109
- odd function, 185
- oddness, 258
- off-diagonal entries, 250
- Old Royal Observatory, 55
- one, 7, 45, 245
- operator, 132
 - + and – as, 133
 - downward multitarget addition, 264
 - elementary, 250
 - general interchange, 261
 - general scaling, 262
 - leftward multitarget addition, 265
 - linear, 134
 - multitarget addition, 263
 - nonlinear, 134
 - quasielementary, 260
 - rightward multitarget addition, 265
 - truncation, 249
 - upward multitarget addition, 264
 - using a variable up, 133
- order, 11
- orientation, 51
- origin, 45, 51
- orthogonal complement, 346
- orthogonal vector, 345
- orthogonalization, 346, 351
- orthonormal rows and columns, 354

- orthonormal vectors, 62
- orthonormalization, 319, 346
- orthonormalizing decomposition, 348
 - inverting a matrix by, 355
- overdetermined linear system, 313, 335, 336
- padding a matrix with zeros, 242
- parallel addition, 118, 220
- parallel subtraction, 121
- parallel unit triangular matrix, 269
 - properties of, 273
- parity, 258, 357
- Parseval's principle, 206
- Parseval, Marc-Antoine (1755–1836), 206
- partial
 - sum, 175
- partial unit triangular matrix, 274
- partial-fraction expansion, 205
- particular solution, 333
- Pascal's triangle, 72
 - neighbors in, 71
- Pascal, Blaise (1623–1662), 72
- path integration, 144
- payment rate, 198
- pencil, 299, 328
- Penrose, Roger (1931–), 335
- permutation, 70
- permutation matrix, 261
- permutor, 261, 357
 - marking, 358
- perpendicular matrix, 346
- perturbation, 206, 383
- Pfufnik, Gorbag, 164
- phase, 40
- physicist, 3
- pilot, 331
- pivot, 289
 - small, 334
- Planck, Max (1858–1947), 29
- plane, 51, 332
- plausible assumption, 112
- point, 51, 62, 332
 - in vector notation, 51
- pole, 38, 84, 162, 164, 171, 187
 - multiple, enclosing a, 173
 - circle of, 206
 - double, 206, 210
 - multiple, 38, 206, 210
 - of a trigonometric function, 186
 - repeated, 206, 210
 - separation of, 206
- polygon, 33
- polynomial, 21, 116
 - characteristic, 368, 382, 384
 - has at least one root, 117
 - of order N has N roots, 117
- positive definiteness, 339, 371
- potentiometer, 331
- power, 15
 - complex, 74
 - fractional, 17
 - integral, 16
 - notation for, 16
 - of a power, 19
 - of a product, 19
 - properties, 16
 - real, 18
 - sum of, 20
- power series, 21, 43
 - bounds on, 175
 - common quotients of, 29
 - derivative of, 75
 - dividing, 22
 - dividing by matching coefficients, 26
 - extending the technique, 29
 - multiplying, 22
 - shifting the expansion point of, 157
 - with negative powers, 21
- prime factorization, 112
- prime mark ($'$), 51
- prime number, 111
 - infinite supply of, 111
 - relative, 231
- processor, computer, 290
- product, 13
 - determinant of, 364
 - dot or inner, 344
 - of determinants, 364
- product rule, derivative, 81, 195, 277

- productivity, 121
- professional mathematician, 114, 212
- professional mathematics, 2, 146
- projector, 390
- proof, 1
 - by contradiction, 112, 307
 - by induction, 41, 154
 - by sketch, 2, 34
- proving backward, 123
- proviso, 326
- pseudoinverse, 335, 352
- pure mathematics, 2, 114
- pyramid
 - volume of, 139
- Pythagoras (c. 580–c. 500 B.C.), 36
- Pythagorean theorem, 36
 - and the hyperbolic functions, 102
 - and the sine and cosine functions, 46
 - in three dimensions, 38
- quadrant, 45
- quadratic expression, 11, 118, 219, 222
- quadratic formula, 12
- quadratics, 11
- quantity
 - exact, 305
 - imprecise, 305
- quartic expression, 11, 118, 219, 227
 - resolvent cubic of, 228
 - roots of, 230
- quartic formula, 230
- quasielementary
 - marking, 358
 - row addition, 264
- quasielementary operator, 260
 - addition, 263
 - interchange, 261
 - scaling, 262
- quintic expression, 11, 118, 231
- quotient, 22, 26, 205
- radian, 45, 55
- radius, 45
- range, 38, 327
- range contour, 163
- rank, 305, 310
 - and independent rows, 296
 - column, 313, 315, 325, 338
 - full, 312
 - impossibility of promoting, 307
 - maximum, 296
 - row, 313, 346
 - uniqueness of, 310
- rank- n determinant, 358
- rank- r inverse, 258, 320
- Raphson, Joseph (1648–1715), 86, 343
- rate, 67
 - relative, 80
- rate of change, 67
- rate of change, instantaneous, 67
- ratio, 18, 205
 - fully reduced, 115
- rational function, 205
 - derivatives of, 210
 - integral of, 209
- rational number, 115
- rational root, 231
- real exponential, 91
- real number, 39
 - approximating as a ratio of integers, 18
- real part, 41
- reciprocal, 7, 323
- reciprocal pair, 322
- rectangle, 7
 - splitting down the diagonal, 34
- rectangular coordinates, 45, 62
- rectangular matrix, 323
- regular part, 172, 187
- relative primeness, 231
- relative rate, 80
- remainder, 22
 - after division by $z - \alpha$, 116
 - zero, 116
- repeated eigensolution, 383
- repeated eigenvalue, 370, 373, 383
- repeated pole, 206, 210
- residual, 334
 - minimizing the, 339

- squared norm of, 335
- residue, 172, 205
- resolvent cubic, 228
- reverse logic, 340
- reversibility, 282, 310
- revolution, 45
- right triangle, 34, 45
- right-hand rule, 51
- rightward multitarget addition operator, 265
- rigor, 2, 136, 156
- rise, 45
- Roman alphabet, 401
- root, 11, 17, 38, 84, 116, 219
 - double, 225
 - finding of numerically, 86, 343
 - guessing, 229
 - rational, 231
 - superfluous, 223
 - triple, 226
- root extraction
 - from a cubic polynomial, 223
 - from a quadratic polynomial, 12
 - from a quartic polynomial, 230
- rotation, 51
 - angle of, 53
- rounding error, 334
- row, 234, 291
 - addition of, 361
 - null, 360
 - null, appending a, 323
 - orthonormal, 354
 - scaled and repeated, 360
- row addition quasialementary, 264
- row operator, 239
- row rank
 - full, 313
- row vector, 344
- Royal Observatory, Old, 55
- rudder, 331
- rudiments, 233, 279
- run, 45
- Saturday, 335
- scalar, 47, 234
 - complex, 50
 - condition of, 375
- scalar matrix, 245
- scaling, 359
- scaling operator
 - elementary, 251
 - general, 262
- scaling quasialementary, 262
- schematic, 299, 328
- Schmidt, Erhard (1876–1959), 346
- Schur decomposition, 377
- Schur, Issai (1875–1941), 266, 377
- screw, 51
- second derivative, 83
- selection from among wooden blocks, 70
- self-adjoint matrix, 384
- separation of poles, 206
- serial addition, 118
- series, 13
 - arithmetic, 15
 - convergence of, 65
 - geometric, 29, 179
 - geometric, variations on, 29
 - harmonic, 179
 - multiplication order of, 14
 - notation for, 13
 - product of, 13
 - sum of, 13
- series addition, 118
- shape
 - area of, 139
- shift operator, 276, 326
- shifting an eigenvalue, 370
- sifting property, 146, 241
- sign
 - alternating, 175
- similarity, 376
- similarity transformation, 254, 282, 375
- Simpson's rule, 129
- simultaneous system of linear equations, 237, 324
- sine, 45
 - derivative of, 103, 106
 - in complex exponential form, 102
 - law of sines, 59

- single-valued function, 162, 164
- singular matrix, 323, 374
 - determinant of, 363
 - versus a nondiagonalizable matrix, 373
- singular value decomposition, 388
- singular value matrix, 388
- singularity, 38
 - essential, 39, 165, 190
- sinusoid, 47
- sixfold integral, 138
- skepticism, 374
- sketch, proof by, 2, 34
- slope, 45, 83
- slow convergence, 181
- slow function, 96
- solid
 - surface area of, 140
 - volume of, 139
- solution, 319
 - error in, 374
 - family of, 332
 - general, 334
 - guessing the form of, 197
 - of least-squares, 337
 - particular, 333
 - particular and homogeneous, 333
- sound, 30
- source, 270
- south, 45, 330
- space, 51, 302, 326
 - addressing, 316
- space and time, 138
- spare column, 295, 315
- sparsity, 246
- sphere, 63
 - surface area of, 140
 - volume of, 142
- spherical coordinates, 62
- split form, 333
- square, 56
 - tilted, 36
- square matrix, 242, 315, 320
 - coincident properties of, 366
 - degenerate, 323
- square root, 18, 39
 - calculation by Newton-Raphson, 88
- square, completing the, 12
- squared residual norm, 335
- squares, least, 335
- squares, sum or difference of, 11
- squaring, 224
- strictly triangular matrix, 266
- strip, tapered, 140
- style, 3, 114, 148
- subtraction
 - parallel, 121
- sum, 13
 - partial, 175
 - weighted, 280
- summation, 13
 - compared to integration, 178
 - convergence of, 65
- superfluous root, 223
- superposition, 109, 172, 360
- surface, 137
- surface area, 140
- surface integral, 138
- surface integration, 140
 - closed, 142
- swimming pool, 332
- symmetry, 185, 304
- symmetry, appeal to, 51
- tall matrix, 315
- tangent, 45
 - derivative of, 106
 - in complex exponential form, 101
- tangent line, 86, 92
- tapered strip, 140
- target, 270
- Tartaglia, Niccolò Fontana (1499–1557), 219
- taxonomy
 - of the linear system, 313
- Taylor expansion, first-order, 74
- Taylor series, 151, 159
 - converting a power series to, 157
 - for specific functions, 174
 - in $1/z$, 190

- integration by, 216
 - multidimensional, 191
 - transposing to a different expansion point, 161
- Taylor, Brook (1685–1731), 151, 159
- technician, 331
- term
 - cross-, 191
 - finite number of, 175
- theory, 324
- time and space, 138
- transformation, linear, 236
- transitivity
 - summational and integrodifferential, 135
- transpose, 240, 361
 - conjugate, 240
 - of a matrix inverse, 257
- trapezoid rule, 129
- triangle, 33, 59
 - area of, 34
 - equilateral, 56
 - right, 34, 45
- triangle inequalities, 34
 - complex, 64
 - vector, 64
- triangular matrix, 266, 382
 - construction of, 267
 - parallel, properties of, 273
 - partial, 274
 - unit parallel, 269
- trigonometric family of functions, 103
- trigonometric function, 45
 - derivative of, 106
 - inverse, 45
 - inverse, derivative of, 108
 - of a double or half angle, 55
 - of a hour angle, 55
 - of a sum or difference of angles, 53
 - poles of, 186
- trigonometrics, complex, 101
- trigonometry, 45
 - properties, 48, 61
- triple integral, 138
- triple root, 226
- triviality, 280
- truncation, 298
 - operator, 249
- tuning, 331
- underdetermined linear system, 313
- uniqueness, 212, 323
 - of matrix rank, 310
- unit, 39, 45, 47
 - imaginary, 39
 - real, 39
- unit basis vector, 47
 - cylindrical, 62
 - spherical, 62
 - variable, 62
- unit circle, 45
- unit lower triangular matrix, 265
- unit magnitude, 374
- unit step function, Heaviside, 145
- unit triangular matrix, 265
 - construction of, 267
 - parallel, properties of, 273
 - partial, 274
- unit upper triangular matrix, 265
- unit vector, 47, 345
- unitary matrix, 354
 - determinant of, 364
- unitary similarity, 376
- unitary transformation, 376
- unity, 7, 45, 47
- unknown coefficient, 197
- unsureness, logical, 114
- up, 45, 332
- upper triangular matrix, 265
- upward multitarget addition operator, 264
- utility variable, 59
- variable, 30
 - assignment, 10
 - change of, 10
 - complex, 5, 78
 - definition notation for, 11
 - dependent, 30, 76
 - independent, 30, 76
 - utility, 59

- variable independent infinitesimal, 130
- variable $d\tau$, 130
- vector, 47, 191, 234
 - angle between two, 345
 - arbitrary, 280, 302
 - building from basis vectors, 375
 - column, 307, 344
 - dot or inner product of two, 344
 - driving, 324, 327
 - elementary, 249, 307, 344
 - elementary, left-over, 309
 - generalized, 191, 234
 - integer, 191
 - magnitude of, 345
 - nonnegative integer, 191
 - normalization of, 346
 - notation for, 47
 - orientation of, 344
 - orthogonal, 345
 - orthogonalization of, 346
 - orthonormal, 62
 - orthonormalization of, 346
 - point, 51
 - replacement of, 302, 330
 - rotation of, 51
 - row, 344
 - row of, 234
 - three-dimensional, 49
 - two-dimensional, 49
 - unit, 47, 345
 - unit basis, 47
 - unit basis, cylindrical, 62
 - unit basis, spherical, 62
 - unit basis, variable, 62
 - zero or null, 280
- vector space, 302, 326, 327
 - addressing, 316, 330
- vector, driving, 333
- vertex, 139
- Vieta's parallel transform, 220
- Vieta's substitution, 220
- Vieta's transform, 220, 221
- Vieta, Franciscus (François Viète, 1540–1603), 219, 220
- visualization, geometrical, 281
- volume, 7, 137, 139
- volume integral, 138
- wave
 - complex, 109
 - propagating, 109
- Weierstrass, Karl Wilhelm Theodor (1815–1897), 155
- weighted sum, 280
- west, 45, 330
- wooden block, 70
- worker, 335
- x86-class computer processor, 290
- zero, 7, 38
 - dividing by, 68
 - matrix, 244
 - padding a matrix with, 242
 - vector, 244
- zero matrix, 244
- zero vector, 244, 280