

TEXTBOOKS in MATHEMATICS

ADVANCED LINEAR ALGEBRA

NICHOLAS LOEHR



CRC Press

Taylor & Francis Group

A CHAPMAN & HALL BOOK

ADVANCED LINEAR ALGEBRA

TEXTBOOKS in MATHEMATICS

Series Editor: Ken Rosen

PUBLISHED TITLES

ABSTRACT ALGEBRA: AN INQUIRY-BASED APPROACH

Jonathan K. Hodge, Steven Schlicker, and Ted Sundstrom

ABSTRACT ALGEBRA: AN INTERACTIVE APPROACH

William Paulsen

ADVANCED CALCULUS: THEORY AND PRACTICE

John Srdjan Petrovic

ADVANCED LINEAR ALGEBRA

Nicholas Loehr

COLLEGE GEOMETRY: A UNIFIED DEVELOPMENT

David C. Kay

COMPLEX VARIABLES: A PHYSICAL APPROACH WITH APPLICATIONS AND MATLAB®

Steven G. Krantz

ESSENTIALS OF TOPOLOGY WITH APPLICATIONS

Steven G. Krantz

INTRODUCTION TO ABSTRACT ALGEBRA

Jonathan D. H. Smith

INTRODUCTION TO MATHEMATICAL PROOFS: A TRANSITION

Charles E. Roberts, Jr.

INTRODUCTION TO PROBABILITY WITH MATHEMATICA®, SECOND EDITION

Kevin J. Hastings

LINEAR ALBEBRA: A FIRST COURSE WITH APPLICATIONS

Larry E. Knop

LINEAR AND NONLINEAR PROGRAMMING WITH MAPLE™: AN INTERACTIVE, APPLICATIONS-BASED APPROACH

Paul E. Fishback

MATHEMATICAL AND EXPERIMENTAL MODELING OF PHYSICAL AND BIOLOGICAL PROCESSES

H. T. Banks and H. T. Tran

ORDINARY DIFFERENTIAL EQUATIONS: APPLICATIONS, MODELS, AND COMPUTING

Charles E. Roberts, Jr.

REAL ANALYSIS AND FOUNDATIONS, THIRD EDITION

Steven G. Krantz

TEXTBOOKS in MATHEMATICS

ADVANCED LINEAR ALGEBRA

NICHOLAS LOEHR

Virginia Polytechnic Institute and State University
Blacksburg, USA



CRC Press

Taylor & Francis Group

Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group an **informa** business
A CHAPMAN & HALL BOOK

CRC Press
Taylor & Francis Group
6000 Broken Sound Parkway NW, Suite 300
Boca Raton, FL 33487-2742

© 2014 by Taylor & Francis Group, LLC
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works
Version Date: 20140306

International Standard Book Number-13: 978-1-4665-5902-8 (eBook - PDF)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Visit the Taylor & Francis Web site at
<http://www.taylorandfrancis.com>

and the CRC Press Web site at
<http://www.crcpress.com>

Dedication

This book is dedicated

to Nanette
and Olivia
and Heather
and Linda
and Zoe.

This page intentionally left blank

Contents

Preface	xvii
I Background on Algebraic Structures	1
1 Overview of Algebraic Systems	3
1.1 Groups	3
1.2 Rings and Fields	4
1.3 Vector Spaces	6
1.4 Subsystems	7
1.5 Product Systems	8
1.6 Quotient Systems	9
1.7 Homomorphisms	11
1.8 Spanning, Linear Independence, Basis, and Dimension	13
1.9 Summary	15
1.10 Exercises	17
2 Permutations	23
2.1 Symmetric Groups	23
2.2 Representing Functions as Directed Graphs	24
2.3 Cycle Decompositions of Permutations	24
2.4 Composition of Cycles	26
2.5 Factorizations of Permutations	27
2.6 Inversions and Sorting	28
2.7 Signs of Permutations	29
2.8 Summary	30
2.9 Exercises	30
3 Polynomials	35
3.1 Intuitive Definition of Polynomials	35
3.2 Algebraic Operations on Polynomials	36
3.3 Formal Power Series and Polynomials	37
3.4 Properties of Degree	39
3.5 Evaluating Polynomials	40
3.6 Polynomial Division with Remainder	41
3.7 Divisibility and Associates	43
3.8 Greatest Common Divisors of Polynomials	44
3.9 GCDs of Lists of Polynomials	45
3.10 Matrix Reduction Algorithm for GCDs	46
3.11 Roots of Polynomials	48
3.12 Irreducible Polynomials	49
3.13 Factorization of Polynomials into Irreducibles	51
3.14 Prime Factorizations and Divisibility	52

3.15 Irreducible Polynomials in $\mathbb{Q}[x]$	53
3.16 Irreducibility in $\mathbb{Q}[x]$ via Reduction Mod p	54
3.17 Eisenstein's Irreducibility Criterion for $\mathbb{Q}[x]$	54
3.18 Kronecker's Algorithm for Factoring in $\mathbb{Q}[x]$	55
3.19 Algebraic Elements and Minimal Polynomials	56
3.20 Multivariable Polynomials	58
3.21 Summary	60
3.22 Exercises	62
II Matrices	71
4 Basic Matrix Operations	73
4.1 Formal Definition of Matrices and Vectors	73
4.2 Vector Spaces of Functions	74
4.3 Matrix Operations via Entries	75
4.4 Properties of Matrix Multiplication	77
4.5 Generalized Associativity	78
4.6 Invertible Matrices	79
4.7 Matrix Operations via Columns	81
4.8 Matrix Operations via Rows	83
4.9 Elementary Operations and Elementary Matrices	84
4.10 Elementary Matrices and Gaussian Elimination	86
4.11 Elementary Matrices and Invertibility	87
4.12 Row Rank and Column Rank	87
4.13 Conditions for Invertibility of a Matrix	89
4.14 Summary	90
4.15 Exercises	92
5 Determinants via Calculations	101
5.1 Matrices with Entries in a Ring	101
5.2 Explicit Definition of the Determinant	102
5.3 Diagonal and Triangular Matrices	103
5.4 Changing Variables	103
5.5 Transposes and Determinants	104
5.6 Multilinearity and the Alternating Property	106
5.7 Elementary Row Operations and Determinants	107
5.8 Determinant Properties Involving Columns	109
5.9 Product Formula via Elementary Matrices	109
5.10 Laplace Expansions	111
5.11 Classical Adjoints and Inverses	113
5.12 Cramer's Rule	114
5.13 Product Formula for Determinants	115
5.14 Cauchy–Binet Formula	116
5.15 Cayley–Hamilton Theorem	118
5.16 Permanents	120
5.17 Summary	121
5.18 Exercises	123

6 Concrete vs. Abstract Linear Algebra	131
6.1 Concrete Column Vectors vs. Abstract Vectors	131
6.2 Examples of Computing Coordinates	133
6.3 Concrete vs. Abstract Vector Space Operations	135
6.4 Matrices vs. Linear Maps	136
6.5 Examples of Matrices Associated with Linear Maps	138
6.6 Vector Operations on Matrices and Linear Maps	140
6.7 Matrix Transpose vs. Dual Maps	141
6.8 Matrix/Vector Multiplication vs. Evaluation of Maps	142
6.9 Matrix Multiplication vs. Composition of Linear Maps	142
6.10 Transition Matrices and Changing Coordinates	143
6.11 Changing Bases	144
6.12 Algebras of Matrices and Linear Operators	145
6.13 Similarity of Matrices and Linear Maps	147
6.14 Diagonalizability and Triangulability	147
6.15 Block-Triangular Matrices and Invariant Subspaces	149
6.16 Block-Diagonal Matrices and Reducing Subspaces	150
6.17 Idempotent Matrices and Projections	151
6.18 Bilinear Maps and Matrices	152
6.19 Congruence of Matrices	153
6.20 Real Inner Product Spaces and Orthogonal Matrices	154
6.21 Complex Inner Product Spaces and Unitary Matrices	155
6.22 Summary	156
6.23 Exercises	161
III Matrices with Special Structure	167
7 Hermitian, Positive Definite, Unitary, and Normal Matrices	169
7.1 Conjugate-Transpose of a Matrix	169
7.2 Hermitian Matrices	171
7.3 Hermitian Decomposition of a Matrix	172
7.4 Positive Definite Matrices	173
7.5 Unitary Matrices	174
7.6 Unitary Similarity	176
7.7 Unitary Triangularization	177
7.8 Simultaneous Triangularization	178
7.9 Normal Matrices and Unitary Diagonalization	180
7.10 Polynomials and Commuting Matrices	181
7.11 Simultaneous Unitary Diagonalization	182
7.12 Polar Decomposition: Invertible Case	183
7.13 Polar Decomposition: General Case	184
7.14 Interlacing Eigenvalues for Hermitian Matrices	185
7.15 Determinant Criterion for Positive Definite Matrices	187
7.16 Summary	188
7.17 Exercises	189
8 Jordan Canonical Forms	195
8.1 Examples of Nilpotent Maps	196
8.2 Partition Diagrams	197
8.3 Partition Diagrams and Nilpotent Maps	198
8.4 Computing Images via Partition Diagrams	199

8.5	Computing Null Spaces via Partition Diagrams	200
8.6	Classification of Nilpotent Maps (Stage 1)	201
8.7	Classification of Nilpotent Maps (Stage 2)	202
8.8	Classification of Nilpotent Maps (Stage 3)	203
8.9	Fitting's Lemma	204
8.10	Existence of Jordan Canonical Forms	205
8.11	Uniqueness of Jordan Canonical Forms	206
8.12	Computing Jordan Canonical Forms	207
8.13	Application to Differential Equations	209
8.14	Minimal Polynomials	210
8.15	Jordan–Chevalley Decomposition of a Linear Operator	211
8.16	Summary	212
8.17	Exercises	213
9	Matrix Factorizations	219
9.1	Approximation by Orthonormal Vectors	220
9.2	Gram–Schmidt Orthonormalization	221
9.3	Gram–Schmidt QR Factorization	222
9.4	Householder Reflections	224
9.5	Householder QR Factorization	226
9.6	LU Factorization	227
9.7	Example of the LU Factorization	229
9.8	LU Factorizations and Gaussian Elimination	230
9.9	Permuted LU Factorizations	232
9.10	Cholesky Factorization	234
9.11	Least Squares Approximation	235
9.12	Singular Value Decomposition	236
9.13	Summary	237
9.14	Exercises	239
10	Iterative Algorithms in Numerical Linear Algebra	245
10.1	Richardson's Algorithm	245
10.2	Jacobi's Algorithm	246
10.3	Gauss–Seidel Algorithm	247
10.4	Vector Norms	248
10.5	Metric Spaces	250
10.6	Convergence of Sequences	250
10.7	Comparable Norms	251
10.8	Matrix Norms	252
10.9	Formulas for Matrix Norms	254
10.10	Matrix Inversion via Geometric Series	255
10.11	Affine Iteration and Richardson's Algorithm	256
10.12	Splitting Matrices and Jacobi's Algorithm	257
10.13	Induced Matrix Norms and the Spectral Radius	257
10.14	Analysis of the Gauss–Seidel Algorithm	259
10.15	Power Method for Finding Eigenvalues	259
10.16	Shifted and Inverse Power Method	261
10.17	Deflation	262
10.18	Summary	262
10.19	Exercises	264

IV The Interplay of Geometry and Linear Algebra	271
11 Affine Geometry and Convexity	273
11.1 Linear Subspaces	273
11.2 Examples of Linear Subspaces	274
11.3 Characterizations of Linear Subspaces	275
11.4 Affine Combinations and Affine Sets	276
11.5 Affine Sets and Linear Subspaces	277
11.6 Affine Span of a Set	278
11.7 Affine Independence	279
11.8 Affine Bases and Barycentric Coordinates	280
11.9 Characterizations of Affine Sets	281
11.10 Affine Maps	282
11.11 Convex Sets	283
11.12 Convex Hulls	284
11.13 Carathéodory's Theorem	284
11.14 Hyperplanes and Half-Spaces in \mathbb{R}^n	286
11.15 Closed Convex Sets	287
11.16 Cones and Convex Cones	289
11.17 Intersection Lemma for V-Cones	290
11.18 All H-Cones Are V-Cones	291
11.19 Projection Lemma for H-Cones	292
11.20 All V-Cones Are H-Cones	293
11.21 Finite Intersections of Closed Half-Spaces	294
11.22 Convex Functions	296
11.23 Derivative Tests for Convex Functions	297
11.24 Summary	298
11.25 Exercises	301
12 Ruler and Compass Constructions	309
12.1 Geometric Constructibility	310
12.2 Arithmetic Constructibility	311
12.3 Preliminaries on Field Extensions	312
12.4 Field-Theoretic Constructibility	314
12.5 Proof that $GC \subseteq AC$	314
12.6 Proof that $AC \subseteq GC$	316
12.7 Algebraic Elements and Minimal Polynomials	320
12.8 Proof that $AC = SQC$	322
12.9 Impossibility of Geometric Construction Problems	323
12.10 Constructibility of the 17-Gon	324
12.11 Overview of Solvability by Radicals	326
12.12 Summary	327
12.13 Exercises	328
13 Dual Spaces and Bilinear Forms	335
13.1 Vector Spaces of Linear Maps	335
13.2 Dual Bases	336
13.3 Zero Sets	337
13.4 Annihilators	338
13.5 Double Dual V^{**}	339
13.6 Correspondence between Subspaces of V and V^*	341

13.7 Dual Maps	342
13.8 Nondegenerate Bilinear Forms	344
13.9 Real Inner Product Spaces	345
13.10 Complex Inner Product Spaces	346
13.11 Comments on Infinite-Dimensional Spaces	348
13.12 Affine Algebraic Geometry	349
13.13 Summary	350
13.14 Exercises	352
14 Metric Spaces and Hilbert Spaces	359
14.1 Metric Spaces	360
14.2 Convergent Sequences	361
14.3 Closed Sets	362
14.4 Open Sets	363
14.5 Continuous Functions	364
14.6 Compact Sets	365
14.7 Completeness	366
14.8 Definition of a Hilbert Space	368
14.9 Examples of Hilbert Spaces	369
14.10 Proof of the Hilbert Space Axioms for $\ell_2(X)$	371
14.11 Basic Properties of Hilbert Spaces	373
14.12 Closed Convex Sets in Hilbert Spaces	374
14.13 Orthogonal Complements	375
14.14 Orthonormal Sets	377
14.15 Maximal Orthonormal Sets	378
14.16 Isomorphism of H and $\ell_2(X)$	379
14.17 Continuous Linear Maps	381
14.18 Dual Space of a Hilbert Space	382
14.19 Adjoint	383
14.20 Summary	384
14.21 Exercises	387
V Modules, Independence, and Classification Theorems	395
15 Finitely Generated Commutative Groups	397
15.1 Commutative Groups	397
15.2 Generating Sets	399
15.3 \mathbb{Z} -Independence and \mathbb{Z} -Bases	400
15.4 Elementary Operations on \mathbb{Z} -Bases	401
15.5 Coordinates and \mathbb{Z} -Linear Maps	402
15.6 UMP for Free Commutative Groups	403
15.7 Quotient Groups of Free Commutative Groups	404
15.8 Subgroups of Free Commutative Groups	405
15.9 \mathbb{Z} -Linear Maps and Integer Matrices	406
15.10 Elementary Operations and Change of Basis	408
15.11 Reduction Theorem for Integer Matrices	411
15.12 Structure of \mathbb{Z} -Linear Maps between Free Groups	413
15.13 Structure of Finitely Generated Commutative Groups	414
15.14 Example of the Reduction Algorithm	415
15.15 Some Special Subgroups	417
15.16 Uniqueness Proof: Free Case	418

15.17 Uniqueness Proof: Prime Power Case	419
15.18 Uniqueness of Elementary Divisors	422
15.19 Uniqueness of Invariant Factors	423
15.20 Uniqueness Proof: General Case	424
15.21 Summary	424
15.22 Exercises	426
16 Axiomatic Approach to Independence, Bases, and Dimension	437
16.1 Axioms	437
16.2 Definitions	438
16.3 Initial Theorems	438
16.4 Consequences of the Exchange Axiom	439
16.5 Main Theorems: Finite-Dimensional Case	440
16.6 Zorn's Lemma	442
16.7 Main Theorems: General Case	443
16.8 Bases of Subspaces	446
16.9 Linear Independence and Linear Bases	446
16.10 Field Extensions	448
16.11 Algebraic Independence and Transcendence Bases	449
16.12 Independence in Graphs	452
16.13 Hereditary Systems	453
16.14 Matroids	454
16.15 Equivalence of Matroid Axioms	455
16.16 Summary	456
16.17 Exercises	457
17 Elements of Module Theory	463
17.1 Module Axioms	463
17.2 Examples of Modules	465
17.3 Submodules	466
17.4 Submodule Generated by a Subset	467
17.5 Direct Products, Direct Sums, and Hom Modules	468
17.6 Quotient Modules	470
17.7 Changing the Ring of Scalars	472
17.8 Fundamental Homomorphism Theorem for Modules	472
17.9 More Module Homomorphism Theorems	474
17.10 Chains of Submodules	476
17.11 Modules of Finite Length	478
17.12 Free Modules	479
17.13 Size of a Basis of a Free Module	481
17.14 Summary	483
17.15 Exercises	486
18 Principal Ideal Domains, Modules over PIDs, and Canonical Forms	493
18.1 Principal Ideal Domains	494
18.2 Divisibility in Commutative Rings	494
18.3 Divisibility and Ideals	495
18.4 Prime and Irreducible Elements	496
18.5 Irreducible Factorizations in PIDs	497
18.6 Free Modules over a PID	498
18.7 Operations on Bases	499

18.8 Matrices of Linear Maps between Free Modules	500
18.9 Reduction Theorem for Matrices over a PID	502
18.10 Structure Theorems for Linear Maps and Modules	504
18.11 Minors and Matrix Invariants	505
18.12 Uniqueness of Smith Normal Form	506
18.13 Torsion Submodules	507
18.14 Uniqueness of Invariant Factors	508
18.15 Uniqueness of Elementary Divisors	509
18.16 $F[x]$ -Module Defined by a Linear Operator	510
18.17 Rational Canonical Form of a Linear Map	512
18.18 Jordan Canonical Form of a Linear Map	513
18.19 Canonical Forms of Matrices	514
18.20 Summary	516
18.21 Exercises	518
VI Universal Mapping Properties and Multilinear Algebra	525
19 Introduction to Universal Mapping Properties	527
19.1 Bases of Free R -Modules	529
19.2 Homomorphisms out of Quotient Modules	529
19.3 Direct Product of Two Modules	531
19.4 Direct Sum of Two Modules	532
19.5 Direct Products of Arbitrary Families of R -Modules	533
19.6 Direct Sums of Arbitrary Families of R -Modules	534
19.7 Solving Universal Mapping Problems	537
19.8 Summary	539
19.9 Exercises	541
20 Universal Mapping Problems in Multilinear Algebra	547
20.1 Multilinear Maps	547
20.2 Alternating Maps	548
20.3 Symmetric Maps	549
20.4 Tensor Product of Modules	550
20.5 Exterior Powers of a Module	553
20.6 Symmetric Powers of a Module	555
20.7 Myths about Tensor Products	557
20.8 Tensor Product Isomorphisms	558
20.9 Associativity of Tensor Products	560
20.10 Tensor Product of Maps	561
20.11 Bases and Multilinear Maps	562
20.12 Bases for Tensor Products of Free R -Modules	564
20.13 Bases and Alternating Maps	565
20.14 Bases for Exterior Powers of Free Modules	566
20.15 Bases for Symmetric Powers of Free Modules	567
20.16 Tensor Product of Matrices	568
20.17 Determinants and Exterior Powers	569
20.18 From Modules to Algebras	571
20.19 Summary	573
20.20 Exercises	577

Appendix: Basic Definitions	583
Sets	583
Functions	583
Relations	584
Partially Ordered Sets	585
Further Reading	587
Bibliography	591

This page intentionally left blank

Preface

What is linear algebra, and how is it used? Upon examining almost any introductory text on linear algebra, we find a standard list of topics that seems to define the subject. On one hand, one part of linear algebra consists of computational techniques for solving linear equations, multiplying and inverting matrices, calculating and interpreting determinants, finding eigenvalues and eigenvectors, and so on. On the other hand, there is a theoretical side to linear algebra involving abstract vector spaces, subspaces, linear independence, spanning sets, bases, dimension, and linear transformations.

But there is much more to linear algebra than just vector spaces, matrices, and linear equations! The goal of this book is to explore a variety of advanced topics in linear algebra, which highlight the rich interconnections linking this subject to geometry, algebra, analysis, combinatorics, numerical computation, and many other areas of mathematics. The book consists of twenty chapters, grouped into six main subject areas (algebraic structures, matrices, structured matrices, geometric aspects of linear algebra, modules, and multilinear algebra). Some chapters approach introductory material from a more sophisticated or abstract viewpoint; other chapters provide elementary expositions of more theoretical concepts; yet other chapters offer unusual perspectives or novel treatments of standard results.

Unlike some advanced mathematical texts, this book has been carefully designed to minimize the dependence of each chapter on material found in earlier chapters. Each chapter has been conceived as a “mathematical vignette” devoted to the development of one specific topic. If you need to learn about Jordan canonical forms, or ruler and compass constructions, or the singular value decomposition, or Hilbert spaces, or QR factorizations, or convexity, or normal matrices, or modules, you may turn immediately to the relevant chapter without first wading through ten chapters of algebraic background or worrying that you will need a theorem covered two hundred pages earlier. We do assume the reader has already encountered the basic linear algebra concepts described earlier (solving linear systems, computing with matrices and determinants, knowing elementary facts about vector spaces and linear maps). These topics are all revisited in a more sophisticated setting at various points throughout the book, but this is not the book you should read to learn the mechanics of Gaussian elimination or matrix multiplication for the first time! Chapter 1 provides a condensed review of some pertinent definitions from abstract algebra. But the vast majority of the book requires very little knowledge of abstract algebra beyond the definitions of fields, vector spaces over a field, subspaces, linear transformations, linear independence, and bases. The last three chapters build on the material on modules covered in Chapter 17, and in a few places one needs to know about permutations (covered in Chapter 2) and polynomials (discussed in Chapter 3).

Although this book focuses on theoretical aspects of linear algebra, giving complete proofs of all results, we supplement and explain the general theory with many specific examples and concrete computations. The level of abstraction gradually increases as one proceeds through the book, as we move from matrices to vector spaces to modules. Except in Chapter 16, we have deliberately avoided presenting the mathematical content as a dry skeleton of itemized definitions, theorems, and proofs. Instead, the material is presented in

a narrative format, providing motivation, examples, and informal discussions to help the reader understand the significance of the theorems and the intuition behind the proofs. Each chapter ends with a summary containing a structured list of the principal definitions and results covered in the chapter, followed by a set of exercises. Most exercises are not used in the main text, and consist of examples, computations, and proofs that will aid the reader in assimilating the material from the chapter. The reader is encouraged to use a computer algebra system to help solve computationally intensive exercises.

This text is designed for use by advanced undergraduates or beginning graduate students, or as a reference. For a second course in linear algebra, one could cover most of the chapters in Parts II, III, and IV, reviewing Part I as needed. For a second course in abstract algebra focusing on the role of linear-algebraic ideas, one could cover Parts I, V, and VI along with selections from the rest of the book (such as Chapter 12 or Chapter 13). The next section describes the contents of each chapter; for more details on what is covered in a given chapter, consult the summary for that chapter.

Synopsis of Topics Covered

Part I: Background on Algebraic Structures.

Chapter 1: Overview of Algebraic Systems. This chapter contains a condensed reference for the abstract algebra background needed in some parts of the book. We review the definitions of groups, commutative groups, rings, integral domains, fields, vector spaces, algebras, and homomorphisms of these structures. We also quickly cover the constructions of subgroups, product groups, quotient groups, and their analogues for other algebraic systems. Finally, we recall some fundamental definitions and results from elementary linear algebra involving linear independence, spanning sets, bases, and dimension. Most of the later chapters in this book require only a small subset of the material covered here, e.g., the definitions of a field, vector space, linear map, subspace, linearly independent list, and ordered basis.

Chapter 2: Permutations. This chapter provides the basic definitions and facts about permutations that are required for the study of determinants and multilinear algebra. After describing the symmetric groups S_n , we show how to visualize functions on a finite set using directed graphs. This idea leads naturally to the factorization of permutations into disjoint cycles. One can also write permutations as products of transpositions or basic transpositions. To understand the role of basic transpositions, we introduce inversions and prove that a list of numbers $f = [f(1), f(2), \dots, f(n)]$ can be sorted into increasing order by interchanging adjacent elements $\text{inv}(f)$ times. For permutations f , we define $\text{sgn}(f) = (-1)^{\text{inv}(f)}$ and use this formula to establish fundamental properties of the sgn function.

Chapter 3: Polynomials. This chapter covers background on one-variable polynomials over a field, which is used in chapters discussing canonical forms, ruler and compass constructions, primary decompositions, commuting matrices, etc. Topics include intuitive and formal definitions of polynomials, evaluation homomorphisms, the universal mapping property for polynomial rings, polynomial division with remainder, greatest common divisors, roots of polynomials, irreducible polynomials, existence and uniqueness of prime factorizations, minimal polynomials, and ways to test polynomials for irreducibility.

Part II: Matrices.

Chapter 4: Basic Matrix Operations. This chapter revisits some basic material about matrices from a somewhat more sophisticated perspective. Topics studied include formal definitions of matrices and matrix operations, vector spaces of functions and matrices, matrix multiplication and its properties, invertible matrices, elementary row and column operations, elementary matrices, Smith canonical form, factorization of invertible matrices into elementary matrices, the equality of row rank and column rank, and the theorem giving equivalent conditions for a matrix to be invertible.

Chapter 5: Determinants via Calculations. The main properties of matrix determinants are often stated without full proofs in a first linear algebra course. This chapter establishes these properties starting from an explicit definition of $\det(A)$ as a sum over permutations of signed products of entries of A . Topics include multilinearity of the determinant as a function of the rows (or columns) of A , the alternating property, the effect of elementary row (or column) operations, the product formula, Laplace expansions, the classical adjoint of A , the explicit formula for A^{-1} , Cramer's rule, the Cauchy–Binet formula, the Cayley–Hamilton theorem, and a quick introduction to permanents.

Chapter 6: Concrete vs. Abstract Linear Algebra. In elementary linear algebra, one first learns to execute computations involving “concrete” objects such as column vectors and matrices. Later, one learns about abstract vector spaces and linear transformations. Concrete concepts defined for matrices (e.g., matrix multiplication, matrix transpose, diagonalizability, idempotence) have abstract counterparts for linear transformations (e.g., composition of maps, dual maps, existence of a basis of eigenvectors, and being a projection). The goal of this chapter is to give a thorough account of the relations between the concrete world of column vectors and matrices on the one hand, and the abstract world of vector spaces and linear maps on the other hand. We build a dictionary linking these two worlds, explaining the exact connection between each abstract concept and its concrete manifestation. In particular, we carefully describe how taking coordinates relative to an ordered basis yields a vector space isomorphism between an abstract vector space V and a concrete space F^n of column vectors; whereas computing the matrix of a linear map relative to an ordered basis gives an algebra isomorphism between an algebra of linear operators and an algebra of matrices. We also discuss how congruence, orthogonal similarity, and unitary similarity of matrices are related to bilinear maps, real inner product spaces, and complex inner product spaces.

Part III: Matrices with Special Structure.

Chapter 7: Hermitian, Positive Definite, Unitary, and Normal Matrices. This chapter develops the basic facts about the special types of matrices mentioned in the title by building an analogy between properties of complex numbers (being real, being positive, having modulus 1) and properties of matrices (being Hermitian, being positive definite, being unitary). This analogy provides motivation for the central notion of a normal matrix. The chapter also includes a discussion of unitary similarity, triangularization, and diagonalization, the spectral theorem for normal matrices, simultaneous triangularization and diagonalization of commuting families, the polar decomposition of a matrix, and the singular value decomposition.

Chapter 8: Jordan Canonical Forms. This chapter gives a novel and elementary derivation of the existence and uniqueness of the Jordan canonical form of a complex matrix. The first step in the proof is to classify nilpotent linear operators (with scalars in any field) by a visual analysis of partition diagrams. Once this classification is complete, we combine it

with a version of Fitting's lemma to derive Jordan canonical forms with no explicit mention of generalized eigenspaces. The chapter concludes with remarks on how one might compute the Jordan form, an application to systems of differential equations, and a discussion of the Jordan–Chevalley decomposition of a linear operator into a sum of commuting nilpotent and diagonalizable linear maps.

Chapter 9: Matrix Factorizations. There are many ways to factor a matrix into products of other matrices that have a special structure (e.g., triangular or unitary matrices). These factorizations often appear as building blocks in algorithms for the numerical solution of linear algebra problems. This chapter proves the existence and uniqueness of several matrix factorizations and explores the algebraic and geometric ideas leading to these factorizations. Topics covered include the Gram–Schmidt orthonormalization algorithm, QR -factorizations, Householder reflections, LU factorizations and their relation to Gaussian elimination, Cholesky's factorization of positive semidefinite matrices, the normal equations for solving least-squares problems, and the singular value decomposition.

Chapter 10: Iterative Algorithms in Numerical Linear Algebra. This chapter studies methods for solving linear systems and computing eigenvalues that employ an iterative process to generate successive approximations converging to the true solution. Iterative algorithms for solving $Ax = b$ include Richardson's method, the Jacobi method, and the Gauss–Seidel method. After developing some background on vector norms and matrix norms, we establish some theoretical results ensuring the convergence of these algorithms for certain classes of matrices. The chapter ends with an analysis of the power method for iteratively computing the largest eigenvalue of a matrix. Techniques for finding other eigenvalues (the shifted power method, inverse power method, and deflation) are also discussed briefly.

Part IV: The Interplay of Geometry and Linear Algebra.

Chapter 11: Affine Geometry and Convexity. Most introductions to linear algebra include an account of vector spaces, linear subspaces, the linear span of a subset of \mathbb{R}^n , linear independence, bases, dimension, and linear transformations. However, the parallel theory of affine sets and affine transformations is seldom covered in detail. This chapter starts by developing the basic notions of affine geometry: affine sets, the affine span of a subset of \mathbb{R}^n , affine combinations, characterizations of affine sets (as translates of linear subspaces, as solution sets of linear equations, and as intersections of hyperplanes), affine independence, barycentric coordinates, and affine transformations. We continue with an introduction to the vast subject of convexity, discussing convex sets, convex hulls, convex combinations, Carathéodory's theorem, simplexes, closed convex sets, convex cones, descriptions of convex sets as intersections of half-spaces, convex functions, epigraphs, Jensen's inequality, and derivative tests for convex functions.

Chapter 12: Ruler and Compass Constructions. Why can a regular 17-gon be constructed with ruler and compass, whereas a regular 7-gon cannot? Which angles can be trisected with a ruler and compass? These questions and others are answered here in a linear-algebraic way by viewing field extensions as vector spaces and looking at dimensions. This material is often presented toward the end of an abstract algebra course, or is skipped altogether. Our development is entirely elementary, using only basic facts about analytic geometry, plane geometry, fields, polynomials, and dimensions of vector spaces.

Chapter 13: Dual Spaces and Bilinear Forms. A fruitful idea in mathematics is the interplay between geometric spaces and structure-preserving functions defined on these spaces. This chapter studies the fundamental case of a finite-dimensional vector space V

and the dual space V^* of linear functions from V to the field of scalars. We show how the notions of *zero sets* and *annihilators* set up inclusion-reversing bijections between the subspaces of V and the subspaces of V^* . Once this correspondence is understood, we explore how the choice of an inner product on V (or, more generally, a nondegenerate bilinear form) sets up an isomorphism between V and V^* . We also discuss the double dual V^{**} and its relation to V , along with dual maps and adjoint operators. The chapter concludes with a few comments on Banach spaces and affine algebraic geometry.

Chapter 14: Metric Spaces and Hilbert Spaces. This chapter gives readers a taste of some aspects of infinite-dimensional linear algebra that play a fundamental role in modern analysis. We begin with a self-contained account of the material we need from analysis, including metric spaces, convergent sequences, closed sets, open sets, continuous functions, compactness, Cauchy sequences, and completeness. Next we develop the basic theory of Hilbert spaces, exploring topics such as the Schwarz inequality, the parallelogram law, orthogonal complements of closed subspaces, orthonormal sets, Bessel's inequality, maximal orthonormal sets, Parseval's equation, abstract Fourier expansions, the role of the Hilbert spaces $\ell^2(X)$, Banach spaces of continuous linear maps, the dual of a Hilbert space, and the adjoint of an operator. A major theme is the interaction between topological properties (especially completeness), algebraic computations, and geometric ideas such as convexity and orthogonality. Beyond a few brief allusions to L^2 spaces, this chapter does not assume any detailed knowledge of measure theory or Lebesgue integration.

Part V: Modules, Independence, and Classification Theorems.

Chapter 15: Finitely Generated Commutative Groups. A central theorem of group theory asserts that every finitely generated commutative group is isomorphic to a product of cyclic groups. This theorem is frequently stated, but not always proved, in first courses on abstract algebra. This chapter proves this fundamental result using linear-algebraic methods. We begin by reviewing basic concepts and definitions for commutative groups, stressing the analogy to corresponding concepts for vector spaces. The key idea in the proof of the classification theorem is to develop the analogue of Gaussian elimination for integer-valued matrices and to use this algorithm to reduce such matrices to a certain canonical form. We also discuss elementary divisors and invariant factors and prove uniqueness results for these objects. The uniqueness proof is facilitated by using partition diagrams to visualize finite groups whose size is a prime power. This chapter uses the very concrete setting of integer matrices to prepare the reader for more abstract algebraic ideas (universal mapping properties and the classification of finitely generated modules over principal ideal domains) covered in later chapters.

Chapter 16: Axiomatic Approach to Independence, Bases, and Dimension. This chapter presents a general axiomatic treatment of some fundamental concepts from linear algebra: linear independence, linear dependence, spanning sets, and bases. A major objective is to prove that every vector space V has a basis, and the cardinality of the basis is uniquely determined by V . One benefit of the axiomatic approach is that it sweeps away a lot of irrelevant extra structure, isolating a few key properties that underlie the main theorems about linear independence and bases. Even better, these axioms arise in other situations besides classical linear algebra. Hence, all the theorems deduced from the axioms will apply to those other situations as well. For example, we will prove properties of the transcendence degree of field extensions by verifying the axioms of the general theory. We conclude with a quick introduction to matroids, which arise by specializing our axiomatic framework to finite sets. This chapter departs from the narrative style of the rest of the book by giving a rigidly structured sequence of axioms, definitions, theorems, and proofs.

Chapter 17: Elements of Module Theory. This chapter introduces the reader to some fundamental concepts in the theory of modules over arbitrary rings. We adopt the approach of using the analogy to vector spaces (which are modules over fields) to motivate fundamental constructions for modules, while carefully pointing out that certain special facts about vector spaces fail to generalize to modules. In particular, the fact that not every module has a basis leads to a discussion of free modules and their properties. Specific topics covered include submodules, quotient modules, direct sums, generating sets, direct products, Hom modules, change of scalars, module homomorphisms, kernels, images, isomorphism theorems, the Jordan–Hölder theorem, length of a module, free modules, bases, and invariance of the size of a basis when the ring of scalars is commutative.

Chapter 18: Principal Ideal Domains, Modules over PIDs, and Canonical Forms. This chapter gives a detailed proof of one of the cornerstones of abstract algebra: the classification of all finitely generated modules over a PID. The chapter begins by proving the necessary algebraic properties of principal ideal domains, including the fact that every PID is a unique factorization domain. The classification proof is modeled upon the concrete case of commutative groups (covered in Chapter 15); a central idea is to develop a matrix reduction algorithm for matrices with entries in a PID. This algorithm changes any matrix into a diagonal matrix called a Smith normal form, which leads to the main structural results for modules. As special cases of the general theory, we deduce theorems on the rational canonical form and Jordan canonical form of linear operators and matrices. We also prove uniqueness of the elementary divisors and invariant factors of a module, as well as uniqueness of the various canonical forms for matrices.

Part VI: Universal Mapping Properties and Multilinear Algebra.

Chapter 19: Introduction to Universal Mapping Properties. The concept of a universal mapping property (UMP) pervades linear and abstract algebra, but is seldom mentioned in introductory treatments of these subjects. This chapter gives a careful and detailed introduction to this idea, starting with the UMP satisfied by a basis of a vector space. The UMP is formulated in several equivalent ways (as a diagram completion property, as a unique factorization property, and as a bijection between collections of functions). The concept is further developed by describing the UMP’s characterizing free R -modules, quotient modules, direct products, and direct sums. This chapter serves as preparation for the next chapter on multilinear algebra.

Chapter 20: Universal Mapping Properties in Multilinear Algebra. The final chapter uses the idea of a universal mapping property (UMP) to organize the development of basic constructions in multilinear algebra. Topics covered include multilinear maps, alternating maps, symmetric maps, tensor products of modules, exterior powers, symmetric powers, and the homomorphisms of these structures induced by linear maps. We also discuss isomorphisms between tensor product modules, bases of tensor products and related modules, tensor products of matrices, the connection between exterior powers and determinants, and tensor algebras.

I wish you an exciting journey through the rich landscape of linear algebra!

Nicholas A. Loehr

Part I

Background on Algebraic Structures

This page intentionally left blank

Overview of Algebraic Systems

This chapter gives a rapid overview of the algebraic systems (such as groups, rings, fields, vector spaces, and algebras) that appear later in the book. After giving the axioms defining each of these systems and some basic examples, we describe some constructions (such as subspaces, product spaces, and quotient spaces) for building new systems from old ones. Then we discuss homomorphisms, which are structure-preserving maps between algebraic systems. The chapter concludes with a review of linear independence, spanning, basis, and dimension in the context of vector spaces over a field.

Unlike the rest of the book, this chapter is intended to be used as a quick review (for readers familiar with abstract algebra) or as a reference (for readers unfamiliar with the definitions), not as a leisurely introduction to the subject. To read the majority of the book, it suffices to know the meaning of the following terms defined in this chapter: commutative ring, field, vector space over a field, linear transformation, subspace, linearly independent list, basis, and dimension. Some further basic definitions regarding sets, functions, relations, and partially ordered sets appear in the Appendix.

1.1 Groups

Given any set S , a *binary operation* on S is a function $p : S \times S \rightarrow S$. We say that S is *closed under* this binary operation to emphasize that $p(a, b)$ is required to belong to S for all $a, b \in S$. There are many operation symbols that are used instead of the notation $p(a, b)$: for example, any of the expressions $a + b$, $a \cdot b$, $a \circ b$, $a \times b$, ab , or $[a, b]$ may be used to abbreviate $p(a, b)$ in different situations.

Next we define some special properties that a binary operation p on S may or may not have. First, p is *commutative* iff¹ $p(a, b) = p(b, a)$ for all $a, b \in S$. Second, p is *associative* iff $p(a, p(b, c)) = p(p(a, b), c)$ for all $a, b, c \in S$. Third, S has an *identity element* relative to p iff there exists $e \in S$ (necessarily unique) such that $p(a, e) = a = p(e, a)$ for all $a \in S$. If such an identity e exists, we say $a \in S$ is *invertible* relative to p iff there exists $a' \in S$ with $p(a, a') = e = p(a', a)$. The element a' is called an *inverse* of a relative to p .

A *group* is a pair (G, p) , where G is a set and p is an associative binary operation on G such that G has an identity element relative to p , and every element of G is invertible relative to p . Writing $p(a, b) = a \star b$, the group axioms take the form shown in Table 1.1.

For example, the set of all nonzero real numbers is a group under multiplication with identity $e = 1$. More generally, the set $\text{GL}_n(\mathbb{R})$ of all $n \times n$ real-valued matrices A having nonzero determinant is a group under matrix multiplication, where the identity element is the $n \times n$ identity matrix. (To check the inverse axiom, one needs the theorem from elementary linear algebra that says A is invertible iff $\det(A) \neq 0$. We prove a more general result in Chapter 5.) For another example, let X be any set, and let $S(X)$ be the set of

¹Throughout this text, the word *iff* is defined to mean “if and only if.”

TABLE 1.1Group Axioms for (G, \star) .

- | |
|---|
| 1. For all $a, b \in G$, $a \star b$ lies in G (closure).
2. For all $a, b, c \in G$, $a \star (b \star c) = (a \star b) \star c$ (associativity).
3. There exists $e \in G$ such that for all $a \in G$, $a \star e = a = e \star a$ (identity).
4. For all $a \in G$, there exists $a^{-1} \in G$ with $a \star a^{-1} = e = a^{-1} \star a$ (inverses). |
|---|

all bijections (one-to-one, onto functions) $f : X \rightarrow X$. Taking the binary operation to be composition of functions, one can verify that $(S(X), \circ)$ is a group. (This group is discussed further in Chapter 2.) The group operations in $\text{GL}_n(\mathbb{R})$ and $S(X)$ are not commutative in general.

A *commutative group* (also called an *Abelian group*) is a group G in which the group operation does satisfy commutativity. Writing $p(a, b) = a + b$ for the group operation, the axioms for a commutative group (in additive notation) take the form shown in Table 1.2.

TABLE 1.2Axioms for a Commutative Group $(G, +)$.

- | |
|---|
| 1. For all $a, b \in G$, $a + b$ lies in G (closure).
2. For all $a, b, c \in G$, $a + (b + c) = (a + b) + c$ (associativity).
3. There exists $0_G \in G$ such that for all $a \in G$, $a + 0_G = a = 0_G + a$ (identity).
4. For all $a \in G$, there exists $-a \in G$ with $a + (-a) = 0_G = (-a) + a$ (inverses).
5. For all $a, b \in G$, $a + b = b + a$ (commutativity). |
|---|

The familiar number systems (the integers \mathbb{Z} , the rational numbers \mathbb{Q} , the real numbers \mathbb{R} , and the complex numbers \mathbb{C}) are all commutative groups under addition. The set \mathbb{R}^k of k -dimensional real vectors $\mathbf{v} = (v_1, \dots, v_k)$ is a commutative group under vector addition; here, $(v_1, \dots, v_k) + (w_1, \dots, w_k) = (v_1 + w_1, \dots, v_k + w_k)$ for all $v_i, w_i \in \mathbb{R}$. The identity element of this group is $\mathbf{0} = (0, \dots, 0)$, and the inverse of $\mathbf{v} = (v_1, \dots, v_k)$ is $-\mathbf{v} = (-v_1, \dots, -v_k)$.

All of the commutative groups just mentioned are infinite. To give examples of finite groups, fix a positive integer n . Let $\mathbb{Z}_n = \{0, 1, 2, \dots, n-1\}$ be the set of *integers modulo n*. Define *addition modulo n* as follows: given $a, b \in \mathbb{Z}_n$, let $a \oplus b = a + b$ if $a + b < n$, and let $a \oplus b = a + b - n$ if $a + b \geq n$. Equivalently, $a \oplus b$ is the remainder when $a + b$ is divided by n , denoted $(a + b) \bmod n$. One may readily verify that (\mathbb{Z}_n, \oplus) is a commutative group of size n .

1.2 Rings and Fields

A *ring* is a triple (R, p, q) , where R is a set and p and q are binary operations on R (denoted by $p(a, b) = a + b$ and $q(a, b) = a \cdot b$) such that: $(R, +)$ is a commutative group; the multiplication operation q is associative with an identity element 1_R ; and the two *distributive*

laws $a \cdot (b + c) = (a \cdot b) + (a \cdot c)$ and $(a + b) \cdot c = (a \cdot c) + (b \cdot c)$ hold for all $a, b, c \in R$. The ring axioms are written out in detail in Table 1.3. Note that some authors do not require, as we do, that all rings have a multiplicative identity.

TABLE 1.3

Ring Axioms for $(R, +, \cdot)$.

1. $(R, +)$ is a commutative group. (See Table 1.2.)
2. For all $a, b \in R$, $a \cdot b$ lies in R (closure under multiplication).
3. For all $a, b, c \in R$, $a \cdot (b \cdot c) = (a \cdot b) \cdot c$ (associativity of multiplication).
4. There is $1_R \in R$ so that for all $a \in R$, $a \cdot 1_R = a = 1_R \cdot a$ (multiplicative identity).
5. For all $a, b, c \in R$, $a \cdot (b + c) = (a \cdot b) + (a \cdot c)$ (left distributive law).
6. For all $a, b, c \in R$, $(a + b) \cdot c = (a \cdot c) + (b \cdot c)$ (right distributive law).

We say R is a *commutative ring* iff its multiplication operation is commutative (for all $a, b \in R$, $a \cdot b = b \cdot a$). We say R is an *integral domain* iff R is a commutative ring such that $0_R \neq 1_R$ and for all $a, b \in R$, if $a \neq 0_R$ and $b \neq 0_R$, then $a \cdot b \neq 0_R$. The last condition states that R has no *zero divisors*, which are nonzero ring elements whose product is zero. We say R is a *field* iff R is a commutative ring such that $0_R \neq 1_R$ and every nonzero element of R has a multiplicative inverse. When R is a field, the set of nonzero elements of R forms a commutative group under multiplication with identity 1_R . One may show that *every field is an integral domain*.

The number systems \mathbb{Z} , \mathbb{Q} , \mathbb{R} , and \mathbb{C} are all commutative rings under addition and multiplication. \mathbb{Q} , \mathbb{R} , and \mathbb{C} are fields, but \mathbb{Z} is not, since most integers do not have multiplicative inverses *within the set* \mathbb{Z} of integers. However, \mathbb{Z} is an integral domain, since the product of two nonzero integers is never zero. For $n \geq 1$, we can make the set \mathbb{Z}_n of integers modulo n into a commutative ring of size n . Addition here is addition mod n , and multiplication is the operation \otimes defined (for $a, b \in \mathbb{Z}_n$) by letting $a \otimes b$ be the remainder when ab is divided by n , denoted $ab \text{ mod } n$. One can verify that the ring $(\mathbb{Z}_n, \oplus, \otimes)$ is a field iff n is a prime number (an integer larger than 1 that is divisible only by itself and 1).

Another example of a commutative ring is the set $\mathbb{R}[x]$ of polynomials in one variable with real coefficients, with the ordinary rules for adding and multiplying polynomials. This ring is not a field, since the only polynomials whose inverses are also polynomials are the nonzero constants. This example can be generalized in several ways. One can replace the real coefficients with coefficients in an arbitrary field (or even an arbitrary ring). Or, one can allow polynomials in more than one variable. For more on polynomials, read Chapter 3.

An example of a non-commutative ring is provided by the set of all $n \times n$ real matrices, under matrix addition and matrix multiplication, where $n \geq 2$ is a fixed integer. More generally, for any ring R , we can consider the set $M_n(R)$ of all $n \times n$ matrices with entries in R . Using the same rules for adding and multiplying matrices as in the real case, this set of matrices becomes a ring that is almost never commutative. For more details, read Chapter 4.

1.3 Vector Spaces

Most introductions to linear algebra study real vector spaces, where the vectors can be multiplied by real numbers called *scalars*. For more advanced investigations, it is helpful to replace the real scalars with elements of more general number systems (such as rings or fields). To ensure that the nice properties of real vector spaces carry over to the more general setting, it turns out that we need to have multiplicative inverses for all nonzero scalars. This leads to the following definition of a *vector space over a field*, which is a direct generalization of the axioms for real vector spaces.

Let F be a field. A *vector space over F* (also called an F -*vector space*) is a triple $(V, +, s)$, where $(V, +)$ is a commutative group, and $s : F \times V \rightarrow V$ is a *scalar multiplication operation* (denoted by $s(c, v) = cv$) satisfying the axioms listed in Table 1.4. If we replace the field F by any ring R , we obtain the definition of an R -*module*. Modules are studied later in the book, starting in Chapter 17. In the context of vector spaces (resp. modules), elements of the field F (resp. the ring R) are often called *scalars*.

TABLE 1.4

Axioms for a Vector Space $(V, +, s)$ over a Field F .

- | |
|---|
| <ol style="list-style-type: none"> 1. $(V, +)$ is a commutative group. (See Table 1.2.) 2. For all $c \in F$ and $v \in V$, $s(c, v) = cv$ lies in V (closure under scalar multiplication). 3. For all $c, d \in F$ and $v \in V$, $(c + d)v = cv + dv$ (distributivity of scalar addition). 4. For all $c \in F$ and $v, w \in V$, $c(v + w) = cv + cw$ (distributivity of vector addition). 5. For all $c, d \in F$ and $v \in V$, $c(dv) = (c \cdot d)v$ (associativity of scalar multiplication). 6. For all $v \in V$, $1_F v = v$ (identity for scalar multiplication). |
|---|

The most well-known example of a real vector space is the set \mathbb{R}^k of vectors $\mathbf{v} = (v_1, \dots, v_k)$ with each $v_i \in \mathbb{R}$. As discussed earlier, \mathbb{R}^k is a commutative group under vector addition. The scalar multiplication operation is given by $c\mathbf{v} = (cv_1, cv_2, \dots, cv_k)$ for $c \in \mathbb{R}$ and $\mathbf{v} \in \mathbb{R}^k$. This example generalizes readily to the case where scalars come from an arbitrary field F . We define F^k to be the set of k -tuples $v = (v_1, v_2, \dots, v_k)$ with each $v_i \in F$. Given such a v and $w = (w_1, w_2, \dots, w_k) \in F^k$ and $c \in F$, define

$$v + w = (v_1 + w_1, v_2 + w_2, \dots, v_k + w_k);$$

$$cv = (c \cdot v_1, c \cdot v_2, \dots, c \cdot v_k).$$

The operations on the right sides of these equations are the given addition and multiplication operations in F . Using these definitions, it is tedious but straightforward to confirm that F^k is an F -vector space.

Some further examples of F -vector spaces, all of which will be discussed in detail later in this book, are: the set $M_{m,n}(F)$ of $m \times n$ matrices with entries in F (Chapter 4); the set of all functions from an arbitrary set X into F (Chapter 4); the set of all linear maps from one F -vector space to another (Chapters 6 and 13); any field K containing F as a subfield (Chapter 12); and the set of polynomials in one or more variables with coefficients in F (Chapter 3).

We conclude with the definition of an *algebra* with coefficients in F , although this concept will only be needed in a few places in the book. For a field F , an F -*algebra* is

a structure $(A, +, \bullet, s)$ such that $(A, +, \bullet)$ is a ring, $(A, +, s)$ is an F -vector space, and the ring multiplication and scalar multiplication are related by the identities

$$c(v \bullet w) = (cv) \bullet w = v \bullet (cw) \quad (c \in F, v, w \in A).$$

(More precisely, A is an *associative F -algebra with identity*.) Some examples of F -algebras are: the set of $n \times n$ matrices with entries in F ; the set of linear maps from a fixed F -vector space V to itself; and the set of polynomials in one or more variables with coefficients in F . Chapter 6 explores the close relationship between the first two of these algebras.

1.4 Subsystems

For each of the algebraic systems mentioned so far, we can construct new algebraic systems of the same kind using a variety of algebraic constructions. Some recurring constructions in abstract algebra are subsystems, direct products, and quotient systems. The next few sections review how these constructions are defined for groups, rings, fields, vector spaces, and algebras. For a more extensive discussion covering the case of modules, read Chapter 17.

Generally speaking, a *subsystem* of a given algebraic system is a subset that is “closed” under the relevant operations. For instance, a *subgroup* of a group (G, \star) is a subset H of G such that: for all $a, b \in H$, $a \star b \in H$; the identity e of G lies in H ; and for all $a \in H$, $a^{-1} \in H$. The first condition says that the subset H of G is *closed under the group operation \star* ; the second condition says that H is *closed with respect to the identity*; and the third condition says that H is *closed under inverses*. When these closure conditions hold, it follows that the set H becomes a group if we restrict the binary operation $\star : G \times G \rightarrow G$ to the domain $H \times H$. (A similar comment applies to the constructions of other types of subsystems below.) H is a *normal subgroup* of G iff H is a subgroup such that $a \star h \star a^{-1} \in H$ for all $a \in G$ and all $h \in H$. The quantity $a \star h \star a^{-1}$ is called a *conjugate of h in G* . So the condition for normality can be stated by saying that H is *closed under conjugation*. Every subgroup of a commutative group is automatically normal.

A *subring* of a ring $(R, +, \cdot)$ is a subset S of R such that $0_R \in S$, $1_R \in S$, and for all $a, b \in S$, $a + b \in S$, $-a \in S$, and $a \cdot b \in S$. In words, S is a subring of R iff S is an additive subgroup of $(R, +)$ that is closed under multiplication and under the multiplicative identity. S is a *left ideal* of R iff S is an additive subgroup such that $a \cdot s \in S$ for all $a \in R$ and $s \in S$; we say that S is *closed under left multiplication by elements of R* . Similarly, S is a *right ideal* of R iff S is an additive subgroup such that $s \cdot a \in S$ for all $a \in R$ and $s \in S$ (i.e., S is *closed under right multiplication by elements of R*). S is an *ideal* of R (also called a *two-sided ideal* for emphasis) iff S is both a left ideal and a right ideal of R . Under our conventions, the different types of ideals in R need not be subrings, since the ideals are not required to contain the multiplicative identity 1_R . In fact, any left, right, or two-sided ideal of R containing 1_R must be the entire ring.

If F is a field, a *subfield* of F is a subring E of F such that for all nonzero $a \in E$, a^{-1} lies in E . So, subfields are required to be closed under addition, subtraction (i.e., additive inverses), multiplication, and division (i.e., multiplicative inverses for nonzero elements), as well as containing the zero and one of the original field.

Now suppose V is a vector space over a field F . W is called a *subspace* of V iff: $0_V \in W$; for all $v, w \in W$, $v + w \in W$; and for all $c \in F$ and $w \in W$, $cw \in W$. In words, a subset W of V is a vector subspace of V iff W is closed under zero, vector addition, and scalar multiplication. The same definition applies if V is a module over a ring R , but now subspaces are called *submodules* or *R -submodules*.

Finally, if A is an algebra over a field F , a subset B of A is a *subalgebra* iff B is both a subring of the ring A and a subspace of the vector space A . So, B contains 0_A and 1_A , and B is closed under addition, additive inverses, ring multiplication, and scalar multiplication.

Here are some examples to illustrate the preceding definitions. For each fixed integer n , the set $n\mathbb{Z} = \{nk : k \in \mathbb{Z}\}$ of integer multiples of n is a subgroup of the additive group $(\mathbb{Z}, +)$, which is a normal subgroup because \mathbb{Z} is commutative. Each set $n\mathbb{Z}$ is also an ideal of the ring \mathbb{Z} , but not a subring except when $n = \pm 1$. One can check, using the division algorithm for integers, that all subgroups of \mathbb{Z} (and hence all ideals of \mathbb{Z}) are of this form. More generally, given any multiplicative group (G, \star) and any fixed $g \in G$, the set $\langle g \rangle = \{g^n : n \in \mathbb{Z}\}$ of powers of g is a subgroup of G , called the *cyclic subgroup generated by g* . This subgroup need not be normal in G . Given a commutative ring R and a fixed $b \in R$, the set $Rb = \{r \cdot b : r \in R\}$ is an ideal of R , called the *principal ideal generated by b* . In general, this does not coincide with $\mathbb{Z}b = \{nb : n \in \mathbb{Z}\}$, which is the additive subgroup of $(R, +)$ generated by b . \mathbb{Z} is a subring of \mathbb{R} that is not a subfield, since \mathbb{Z} is not closed under multiplicative inverses. \mathbb{Q} is a subfield of \mathbb{R} , as is the set $\mathbb{Q}(\sqrt{2})$ of all real numbers of the form $a + b\sqrt{2}$, where $a, b \in \mathbb{Q}$. The set of polynomials of degree at most 3 (together with zero) is a subspace of the vector space of all polynomials with real coefficients. For any field F , the set $\{(t, t+u, -u) : t, u \in F\}$ is a subspace of the F -vector space F^3 . The set of upper-triangular $n \times n$ matrices is a subalgebra of the \mathbb{R} -algebra of real $n \times n$ matrices.

1.5 Product Systems

Next we consider the *direct product* construction for vector spaces. Suppose V_1, V_2, \dots, V_n are given vector spaces over the same field F . The *product set* $V = V_1 \times V_2 \times \dots \times V_n$ consists of all ordered n -tuples $v = (v_1, v_2, \dots, v_n)$ with each $v_i \in V_i$. We can turn this product set into an F -vector space by defining addition and scalar multiplication as follows. Given $v = (v_1, \dots, v_n)$ and $w = (w_1, \dots, w_n)$ in V and given $c \in F$, define $v + w = (v_1 + w_1, \dots, v_n + w_n)$ and $cv = (cv_1, \dots, cv_n)$. In these definitions, the operations in the i 'th component are the sum and scalar multiplication operations in the vector space V_i . It is tedious but routine to check the vector space axioms for V . In particular, the additive identity of V is $0_V = (0_{V_1}, 0_{V_2}, \dots, 0_{V_n})$, and the additive inverse of $v = (v_1, \dots, v_n) \in V$ is $-v = (-v_1, \dots, -v_n)$. The set V with these operations is called the *direct product* of the F -vector spaces V_1, V_2, \dots, V_n . Note that the F -vector space F^n is a special case of this construction obtained by taking every $V_i = F$.

There are analogous constructions for the direct products of groups, rings, modules, and algebras. In each case, all relevant algebraic operations are defined “componentwise.” For example, given positive integers n_1, n_2, \dots, n_k , the direct product $G = \mathbb{Z}_{n_1} \times \mathbb{Z}_{n_2} \times \dots \times \mathbb{Z}_{n_k}$ is a commutative group of size $n_1 n_2 \cdots n_k$, with operation

$$(a_1, a_2, \dots, a_k) + (b_1, b_2, \dots, b_k) = ((a_1 + b_1) \bmod n_1, (a_2 + b_2) \bmod n_2, \dots, (a_k + b_k) \bmod n_k)$$

for $a_i, b_i \in \mathbb{Z}_{n_i}$. In Chapter 15 we will prove that *every* finite commutative group is isomorphic to a direct product of this form. (Isomorphisms are defined in §1.7 below.)

1.6 Quotient Systems

This section describes the notion of a *quotient system* of a given algebraic system. Quotient constructions are more subtle and less intuitive than the constructions of subsystems and product systems, but they play a prominent role in abstract algebra and other parts of mathematics.

The simplest instance of the quotient construction occurs when $(G, +)$ is a commutative group with subgroup H . We will build a new commutative group, denoted G/H and called the *quotient group of G by H* . Intuitively, the group G/H will be a “simplification” of G obtained by “discarding information contained in H .” For instance, if G is the additive group \mathbb{Z} of integers and H is the subgroup $n\mathbb{Z}$ of multiples of n , then $\mathbb{Z}/n\mathbb{Z}$ will turn out to be isomorphic to the group \mathbb{Z}_n of integers mod n . Intuitively, formation of the quotient will discard multiples of n and focus attention on the remainders when various integers are divided by n .

We now give the details of the construction of G/H , where $(G, +)$ is a commutative group with subgroup H . For each $a \in G$, define the *coset of H represented by a* to be $a + H = \{a + h : h \in H\}$, which is a certain subset of G . Define $G/H = \{a + H : a \in G\}$ to be the set of all cosets of H in G . Thus, G/H is the set of subsets of G obtained by “translating” the subset H by various group elements a .

A critical point is the fact that each coset (element of G/H) can have several different *names*. In other words, there are often multiple ways of writing a particular coset in the form $a + H$ for some $a \in G$. The *coset equality theorem* gives a precise criterion for when two elements of G give us different names for the same coset. The theorem states that for all $a, b \in G$, $a + H = b + H$ iff $a - b \in H$. To prove this, first assume $a, b \in G$ satisfy $a + H = b + H$. Now, since $0_G \in H$, we have $a = a + 0_G \in a + H$. Since the set $a + H$ equals the set $b + H$ by assumption, we then have $a \in b + H$, which means $a = b + h$ for some $h \in H$. But then $a - b = a + (-b) = h \in H$, as needed. Conversely, let us now fix $a, b \in G$ satisfying $a - b \in H$; we must prove $a + H = b + H$. First we prove the set inclusion $a + H \subseteq b + H$. Fix $x \in a + H$; then we can write $x = a + h_1$ for some $h_1 \in H$. Observe that $x = a + h_1 = b + (a - b) + h_1 = b + k$, where $k = (a - b) + h_1$ is the sum of two elements of H , hence is in H . So $x \in b + H$. For the opposite set inclusion, fix $y \in b + H$. Write $y = b + h_2$ for some $h_2 \in H$. Then notice that $y = a + [-(a - b) + h_2] \in a + H$, where $a - b \in H$, $-(a - b) \in H$, and $-(a - b) + h_2 \in H$ because H is a subgroup. So $y \in a + H$ as required.

We note next that G is the disjoint union of the distinct cosets of H . On one hand, every coset of H is a subset of G (by the closure axiom of G). On the other hand, for any $a \in G$, a lies in the coset $a + H$ since $a = a + 0_G$ with $0_G \in H$. To see that any two distinct cosets of H must be disjoint, suppose $a + H$ and $b + H$ are two cosets of H such that some element c lies in both cosets (here $a, b, c \in G$). Write $c = a + h = b + k$ for some $h, k \in H$; then $a - b = k - h \in H$, so the coset equality theorem gives $a + H = b + H$.

Here is an example. In the group $(\mathbb{Z}_{12}, \oplus)$, consider the subgroup $H = \{0, 3, 6, 9\}$. The cosets of H are $0 + H = \{0, 3, 6, 9\} = 3 + H = 6 + H = 9 + H$, $1 + H = \{1, 4, 7, 10\} = 4 + H = 7 + H = 10 + H$, and $2 + H = \{2, 5, 8, 11\} = 5 + H = 8 + H = 11 + H$. Note that any element b in a coset $a + H$ can be used to give a new name $b + H$ for that coset. The set $G/H = \{a + H : a \in \mathbb{Z}_{12}\} = \{\{0, 3, 6, 9\}, \{1, 4, 7, 10\}, \{2, 5, 8, 11\}\}$ consists of three distinct cosets, and the set G is the disjoint union of these cosets.

Returning to the general case, the next step is to define a binary addition operation p on G/H by setting $p(a + H, b + H) = (a + b) + H$ for all $a, b \in G$. (If we reuse the plus symbol $+$ to denote the new operation p , the definition reads: $(a + H) + (b + H) = (a + b) + H$.) Now a

new subtlety emerges: our definition of the “sum” of two elements of G/H depends on the particular *names* $a + H$ and $b + H$ that we have chosen to represent these elements. If we change from these names to other names of the same two cosets, how can we be sure that the output is not affected? To answer this question, suppose $a, a_1, b, b_1 \in G$ satisfy $a + H = a_1 + H$ and $b + H = b_1 + H$. Using the original names $a + H$ and $b + H$, the sum of these cosets was defined to be $(a + b) + H$. Using the new names $a_1 + H$ and $b_1 + H$ for these same cosets, the sum of the cosets would be calculated as $(a_1 + b_1) + H$. These two answers might initially appear to be different, but in fact they are just different *names* for the same coset. To see this, note from the coset equality theorem that $a - a_1 \in H$ and $b - b_1 \in H$. Therefore, since the subgroup H is closed under addition, $(a + b) - (a_1 + b_1) = (a - a_1) + (b - b_1) \in H$. (This calculation also requires commutativity of G .) Another application of the coset equality theorem gives $(a + b) + H = (a_1 + b_1) + H$, as needed. We summarize this calculation by saying that the binary operation p is *well-defined* (or *single-valued*). One must perform such a check every time a binary operation or function is defined in terms of the *name* of an object that has several possible names.

With this technical issue out of the way, we can now verify the commutative group axioms for $(G/H, p)$ without difficulty. Fix $a, b, c \in G$. To check closure of G/H under p , note $p(a + H, b + H) = (a + H) + (b + H) = (a + b) + H$ does lie in G/H , since $a + b \in G$ by the closure axiom for G . To check associativity of p , note

$$\begin{aligned} p(a + H, p(b + H, c + H)) &= p(a + H, (b + c) + H) = (a + (b + c)) + H \\ &= ((a + b) + c) + H = p((a + b) + H, c + H) = p(p(a + H, b + H), c + H), \end{aligned}$$

where the third equality holds by associativity of $+$ in G . Similarly, commutativity of p holds since $(a + H) + (b + H) = (a + b) + H = (b + a) + H = (b + H) + (a + H)$, which uses commutativity of $+$ in G . The identity element of G/H is the coset $0 + H = \{0 + h : h \in H\} = H$, because

$$(a + H) + (0 + H) = (a + 0) + H = a + H = (0 + a) + H = (0 + H) + (a + H)$$

by the identity axiom for $(G, +)$. Finally, the inverse of $a + H$ relative to p is $(-a) + H$, since the inverse axiom in $(G, +)$ gives $(a + H) + ((-a) + H) = (a + (-a)) + H = 0 + H = ((-a) + a) + H = ((-a) + H) + (a + H)$.

Continuing our earlier example where $G = \mathbb{Z}_{12}$ and $H = \{0, 3, 6, 9\}$, recall that $G/H = \{0 + H, 1 + H, 2 + H\}$. We have $(1 + H) + (2 + H) = (1 + 2) + H = 3 + H$, and another name for this answer is $0 + H$. Furthermore, $(2 + H) + (2 + H) = 4 + H = 1 + H$. Comparing the addition table for G/H to the addition table for the group (\mathbb{Z}_3, \oplus) , one can check that the groups \mathbb{Z}_{12}/H and \mathbb{Z}_3 are isomorphic (as defined in §1.7). Similarly, it can be shown that for $n \geq 1$, the quotient group $\mathbb{Z}/n\mathbb{Z}$ is isomorphic to the group (\mathbb{Z}_n, \oplus) defined in §1.1.

The quotient construction just given can now be generalized to other types of algebraic structures. First, we can replace the commutative group $(G, +)$ and its subgroup H by an arbitrary group (G, \star) and a *normal* subgroup H . For each $a \in G$, we now define the *left coset* $a \star H = \{a \star h : h \in H\}$, and we let $G/H = \{a \star H : a \in G\}$. In this setting, the *left coset equality theorem* states that for all $a, b \in G$, $a \star H = b \star H$ iff $a^{-1} \star b \in H$ iff $b^{-1} \star a \in H$ (note the inverse occurs to the *left* of the star in each case). The binary operation p on the set G/H of left cosets of H is now defined by $p(a \star H, b \star H) = (a \star b) \star H$ for all $a, b \in G$. The verification that p is well-defined (which we leave as an exercise) depends critically on the normality of the subgroup H . Once this has been done, it is straightforward to prove (as above) that $(G/H, p)$ is a group. A similar construction could be executed using *right cosets* $H \star a = \{h \star a : h \in H\}$, which satisfy $H \star a = H \star b$ iff $a \star b^{-1} \in H$. However, normality of H implies that every left coset $a \star H$ equals the right coset $H \star a$, so that the resulting quotient group is the same as before.

Second, we can replace G and H by a ring $(R, +, \cdot)$ and an ideal I of R . The set of (additive) cosets $R/I = \{a + I : a \in R\}$ is already known to be a commutative group, since $(R, +)$ is a commutative group with subgroup I . We introduce a second binary operation q on R/I by setting $q(a + I, b + I) = (a \cdot b) + I$ for $a, b \in R$. (Denoting q by \bullet , this takes the form $(a + I) \bullet (b + I) = (a \cdot b) + I$.) One may check as an exercise that q is well-defined, using the assumption that I is a (two-sided) ideal. It is then routine to verify that R/I is a ring with additive identity $0_R + I$ and multiplicative identity $1_R + I$, which is commutative if R is commutative. This ring is called the *quotient ring* of R by the ideal I .

Third, we can replace G and H by an F -vector space $(V, +, s)$ and a subspace W , where F is a field. We will construct the *quotient vector space* V/W . As before, the quotient set $V/W = \{v + W : v \in V\}$ is already known to be a commutative group, since W is a subgroup of $(V, +)$. We introduce a new scalar multiplication $t : F \times V/W \rightarrow V/W$ by setting $t(c, v + W) = (s(c, v)) + W$ for all $c \in F$ and $v \in V$. (Writing the old scalar multiplication using juxtaposition and the new scalar multiplication using \cdot , this formula becomes: $c \cdot (v + W) = (cv) + W$.) Let us check that t is well-defined. Elements of F do not have multiple names, but elements of V/W do. So, we fix $c \in F$ and $u, v \in V$ and assume $u + W = v + W$. Is it true that $t(c, u + W) = t(c, v + W)$? In other words, is $(cu) + W = (cv) + W$? To check this, apply the coset equality theorem to our assumption to see that $u - v \in W$. Since W is a vector subspace of V and $c \in F$, we deduce that $c(u - v) = cu - cv$ also lies in W . Then the coset equality theorem gives $(cu) + W = (cv) + W$, as needed. It is now straightforward to verify the F -vector space axioms for $(V/W, +, t)$. We already have seen that $(V/W, +)$ is a commutative group. To verify distributivity of vector addition in this space, fix $c \in F$ and $v, x \in V$, and calculate

$$c \cdot ((v + W) + (x + W)) = c \cdot ((v + x) + W) = (c(v + x)) + W.$$

By the known distributivity of vector addition in the original space V , the calculation continues:

$$(c(v + x)) + W = (cv + cx) + W = (cv + W) + (cx + W) = c \cdot (v + W) + c \cdot (x + W),$$

and the axiom is verified. The remaining axioms are checked similarly. The construction in this paragraph goes through verbatim in the case where R is a ring and V is an R -module with submodule W (see Chapter 17).

1.7 Homomorphisms

The concept of a *homomorphism* allows us to compare algebraic structures. There are many different kinds of homomorphisms, one for each kind of algebraic system. Intuitively, a homomorphism of an algebraic system is a function from one system to another system of the same kind that “preserves” all relevant operations and structure.

For instance, if (G, \star) and $(K, *)$ are groups, a *group homomorphism* from G to K is a function $T : G \rightarrow K$ such that $T(x \star y) = T(x) * T(y)$ for all $x, y \in G$. It follows from this condition that $T(e_G) = e_K$, $T(x^{-1}) = T(x)^{-1}$, and more generally $T(x^n) = T(x)^n$ for all $x \in G$ and all $n \in \mathbb{Z}$. We say that the group homomorphism T *preserves* the group operations, the identity, inverses, and powers of group elements. In the case where G and K are commutative groups with operations written in additive notation, the definition of a group homomorphism becomes $T(x + y) = T(x) + T(y)$ for all $x, y \in G$. Now $T(0_G) = 0_K$, $T(-x) = -T(x)$, and $T(nx) = nT(x)$ for all $x \in G$ and all $n \in \mathbb{Z}$. (The notation nx denotes the sum of n copies of x for $n > 0$, or the sum of $|n|$ copies of $-x$ for $n < 0$.)

Analogously, given two rings R and S , a *ring homomorphism* is a function $T : R \rightarrow S$ such that for all $x, y \in R$, $T(x + y) = T(x) + T(y)$, $T(x \cdot y) = T(x) \cdot T(y)$, and $T(1_R) = 1_S$. It follows from these conditions that $T(x^n) = T(x)^n$ for all $x \in R$ and all integers $n \geq 0$; the formula also holds for negative integers n when x is an invertible element of R .

Next, suppose V and W are vector spaces over a field F . An *F -linear map* (also called a *vector space homomorphism* or *linear transformation*) is a function $T : V \rightarrow W$ such that $T(u + v) = T(u) + T(v)$ and $T(cv) = cT(v)$ for all $u, v \in V$ and all $c \in F$. The same definition applies if V and W are modules over a ring R ; in this case, we call T an *R -linear map* or *R -module homomorphism*. Finally, an *algebra homomorphism* is a map between two F -algebras that is both a ring homomorphism and a vector space homomorphism.

Let $T : X \rightarrow Y$ be a homomorphism of any of the types defined above. We call T an *isomorphism* iff T is a bijective (one-to-one and onto) function from X to Y . In this case, T has a two-sided inverse function $T^{-1} : Y \rightarrow X$. One may check that T^{-1} is always a homomorphism of the same type as T , and is also an isomorphism. Furthermore, the composition of homomorphisms (resp. isomorphisms) is a homomorphism (resp. isomorphism), and the identity map on a given algebraic structure X is an isomorphism for that type of structure. We write $X \cong Y$ if there exists an isomorphism between X and Y . The preceding remarks show that \cong is an equivalence relation on any fixed set of algebraic structures of a given kind.

The constructions in the preceding section furnish some examples of homomorphisms. If H is a subgroup of G , the inclusion map $i : H \rightarrow G$ given by $i(h) = h$ for all $h \in H$ is an injective group homomorphism; i is an isomorphism iff $H = G$. If H is normal in G , the projection map $p : G \rightarrow G/H$, given by $p(x) = x \star H$ for all $x \in G$, is a surjective group homomorphism; p is an isomorphism iff $H = \{e_G\}$. Similarly, if R is a ring with ideal I , the projection $p : R \rightarrow R/I$ given by $p(x) = x+I$ for $x \in R$ is a surjective ring homomorphism. If V is an F -vector space with subspace W , we obtain an analogous projection $p : V \rightarrow V/W$, which is F -linear. In the case of groups, vector spaces, and modules, the injective maps $j_i : V_i \rightarrow V_1 \times \cdots \times V_i \times \cdots \times V_n$ sending $x_i \in V_i$ to $(0, \dots, x_i, \dots, 0)$ are homomorphisms. However, this statement does not hold for products of rings, since $j_i(1_{V_i})$ is usually not the multiplicative identity of the product ring. On the other hand, for all the types of algebraic systems discussed, the projections $q_i : V_1 \times \cdots \times V_i \times \cdots \times V_n \rightarrow V_i$ sending (x_1, \dots, x_n) to x_i are surjective homomorphisms.

We now discuss kernels, images, and the fundamental homomorphism theorem. Suppose $T : X \rightarrow Y$ is a function mapping a set X into an additive group Y with identity element 0_Y . The *kernel* of T is $\ker(T) = \{x \in X : T(x) = 0_Y\} \subseteq X$, and the *image* of T is $\text{img}(T) = \{T(x) : x \in X\} \subseteq Y$. If T is a group homomorphism, one checks that $\ker(T)$ is a normal subgroup of X and $\text{img}(T)$ is a subgroup of Y . If T is a ring homomorphism, then $\ker(T)$ is an ideal of X and $\text{img}(T)$ is a subring of Y . If T is a linear transformation of F -vector spaces, then $\ker(T)$ and $\text{img}(T)$ are subspaces of X and Y , respectively. If T is an R -module homomorphism, then $\ker(T)$ and $\text{img}(T)$ are R -submodules of X and Y , respectively. In all these cases, T is injective iff $\ker(T) = \{0_X\}$, while T is surjective iff $\text{img}(T) = Y$. When T is a linear map between vector spaces, $\ker(T)$ is sometimes called the *null space* of T , and $\text{img}(T)$ is also called the *range* of T .

Let X and Y be commutative groups with operations written in additive notation. Suppose $T : X \rightarrow Y$ is a group homomorphism with kernel K and image I . The *fundamental homomorphism theorem for groups* asserts that there is a group isomorphism $T' : X/K \rightarrow I$ given by $T'(x+K) = T(x)$ for all $x \in X$. We prove this as follows. (a) Is T' well-defined? Fix $w, x \in X$ and assume $x+K = w+K$ in X/K ; we must check that $T'(x+K) = T'(w+K)$, i.e., that $T(x) = T(w)$. Now, the coset equality theorem gives $(-x)+w \in K = \ker(T)$, so $T((-x)+w) = 0_Y$, so $-T(x)+T(w) = 0_Y$, so $T(x) = T(w)$. (b) Is T' a group

homomorphism? Fix $u, x \in X$, and calculate

$$T'((u+K)+(x+K)) = T'((u+x)+K) = T(u+x) = T(u)+T(x) = T'(u+K)+T'(x+K).$$

(c) Is T' surjective? Note first that T' does map into the codomain I , since $T'(x+K) = T(x) \in I = \text{img}(T)$ for all $x \in X$. Conversely, given any $y \in \text{img}(T)$, the definition of image shows that $y = T(x)$ for some $x \in X$, hence $y = T'(x+K)$ for some $x+K \in X/K$.
(d) Is T' one-to-one? Fix $u, x \in X$ satisfying $T'(u+K) = T'(x+K)$; we must prove $u+K = x+K$. We know $T(u) = T(x)$, so $-T(x) + T(u) = 0_Y$, so $T((-x)+u) = 0_Y$, so $-x+u \in \ker(T) = K$. By the coset equality theorem, $u+K = x+K$ as needed. The proof is now complete; we note that, except for a notation change, the same proof works for arbitrary groups.

There are analogous fundamental homomorphism theorems for rings, vector spaces, modules, and algebras. For example, if X and Y are rings and $T : X \rightarrow Y$ is a ring homomorphism with kernel K and image I , then there exists a ring isomorphism $T' : X/K \rightarrow I$ given by $T'(x+K) = T(x)$ for all $x \in X$. To prove this, note T is a group homomorphism between the additive groups X and Y . By the previous proof, we already know that T' as given in the theorem statement is a bijective, well-defined homomorphism of additive groups. We need only check that T' also preserves the ring multiplication and identity. Fix $u, x \in X$, and calculate

$$T'((u+K) \cdot (x+K)) = T'((u \cdot x)+K) = T(u \cdot x) = T(u) \cdot T(x) = T'(u+K) \cdot T'(x+K).$$

Moreover, $T'(1_{X/K}) = T'(1_X+K) = T(1_X) = 1_Y = 1_I$, completing the proof. We let the reader formulate and prove the analogous homomorphism theorems for vector spaces, modules, and algebras (cf. Chapter 17, which discusses this theorem and other isomorphism theorems in the context of modules).

1.8 Spanning, Linear Independence, Basis, and Dimension

Introductions to linear algebra often discuss the concepts of linear independence, spanning sets, bases, and dimension in the case of real vector spaces. The same ideas occur in the more general setting of vector spaces over a field. In this section, after reviewing the basic definitions, we state without proof some fundamental theorems concerning linear independence and bases of vector spaces. These theorems will be proved, in an even more general context, in Chapter 16.

Let V be a vector space over a field F . Let $L = (v_1, v_2, \dots, v_k)$ be a finite list of vectors with each $v_i \in V$. Any expression of the form $c_1v_1 + c_2v_2 + \dots + c_kv_k$ with each $c_i \in F$ is called an *F -linear combination* of the v_i 's. We say that the list L spans V iff every $v \in V$ can be written in *at least one* way as an F -linear combination of v_1, \dots, v_k . We say that the list L is *linearly dependent over F* (or *F -linearly dependent*) iff there is a list of scalars $(c_1, c_2, \dots, c_k) \neq (0, 0, \dots, 0)$ with $0 = c_1v_1 + c_2v_2 + \dots + c_kv_k$. In other words, L is linearly dependent provided that zero can be written as a linear combination of the vectors in L where at least one coefficient is not the zero scalar. L is *linearly independent over F* iff L is not linearly dependent over F . Spelling this out, L is F -linearly independent iff for all $c_1, \dots, c_k \in F$, $c_1v_1 + \dots + c_kv_k = 0$ implies $c_1 = c_2 = \dots = c_k = 0$.

Linear independence can be rephrased in a way that looks more like the definition of spanning. We claim L is F -linearly independent iff every $v \in V$ can be written in *at most one* way as an F -linear combination of the vectors in L . On one hand, if L is linearly dependent,

then the vector $v = 0$ can be written in at least two ways as a linear combination of the v_i 's: one way is $0 = 0v_1 + 0v_2 + \cdots + 0v_k$, and the other way is the linear combination appearing in the definition of linear dependence. On the other hand, assume L is linearly independent. Suppose a given $v \in V$ can be written as $v = \sum_{i=1}^k a_i v_i$ and also as $v = \sum_{i=1}^k b_i v_i$ with $a_i, b_i \in F$. Subtracting these equations gives $0 = \sum_{i=1}^k (a_i - b_i) v_i$. The assumed linear independence of the v_i 's then gives $a_i - b_i = 0$ for all i , hence $a_i = b_i$ for $1 \leq i \leq k$. So there is at most one way of writing v as a linear combination of the v_i 's.

Continuing the definitions, we say that the list L is an *ordered basis* of V iff L spans V and L is F -linearly independent. This means that every $v \in V$ can be written in *exactly one way* as an F -linear combination of vectors in L . The vector space V is called *finite-dimensional* iff V has a finite ordered basis. A theorem (proved in Chapter 16) assures us that any two ordered bases of a finite-dimensional vector space V must have the same size, where the size of a list is the number of vectors in it. The *dimension* of a finite-dimensional vector space V is the common size of all ordered bases of V .

For vector spaces that might be infinite-dimensional, it is convenient to augment the definitions of spanning, linear independence, and bases to apply to *sets* of vectors as well as *lists* of vectors. Let S be any subset of vectors in an F -vector space V ; S could be finite or infinite. We say S *spans* V iff every $v \in V$ can be written as an F -linear combination of a finite list of vectors (v_1, \dots, v_n) with each $v_i \in S$ (the list used depends on v). We say S is *linearly independent* iff every finite list of distinct elements of S is linearly independent (in the old sense). We say S is a *basis* of V iff S spans V and is linearly independent. As before, S is a basis iff every $v \in V$ can be written uniquely as a finite linear combination of elements of S with nonzero coefficients.

Here are some theorems about these concepts, to be proved in Chapter 16. (Some of these results depend on the axiom of choice from set theory.) Assume $S \subseteq T \subseteq V$, where V is an F -vector space. If S spans V , then T spans V . If T is linearly independent, then S is linearly independent. V always spans V ; the empty set \emptyset is always linearly independent. If T spans V , there exists a basis B of V contained in T ; in particular, letting $T = V$ shows that V has at least one basis. If S is linearly independent, there exists a basis C of V containing S ; in particular, letting $S = \emptyset$ shows again that V has at least one basis. If S is any linearly independent subset of V and U is any spanning set for V , then $|S| \leq |U|$; i.e., there exists a one-to-one map $g : S \rightarrow U$. Any two bases of V have the same cardinality; i.e., if B and C are both bases of V , there exists a bijection $f : B \rightarrow C$. Because of these theorems, we can define the *dimension* of V , denoted $\dim(V)$ or $\dim_F(V)$, to be the unique cardinality of any basis of V . Recall that linearly independent subsets of V cannot have larger cardinality than spanning subsets of V . It follows that any set $S \subseteq V$ with $|S| > \dim(V)$ must be linearly dependent over F , and any set $T \subseteq V$ with $|T| < \dim(V)$ cannot span V .

We should also note that a list $L = (v_1, \dots, v_k)$ spans V iff the set $\{v_1, \dots, v_k\}$ spans V , whereas L is linearly independent iff L contains no repetitions and the set $\{v_1, \dots, v_k\}$ is linearly independent. It follows that, in the finite-dimensional case, any ordered basis of V has the same size as any basis of V , so that the two definitions of dimension are consistent.

As an application of the results stated above, let us prove the *rank-nullity theorem*: given finite-dimensional F -vector spaces V and W and a linear map $T : V \rightarrow W$ with null space $N = \ker(T)$ and range $R = \text{img}(T)$, we have $\dim(N) + \dim(R) = \dim(V)$. To prove this, start with an ordered basis $B_N = (v_1, \dots, v_k)$ of the subspace N (so $\dim(N) = k$). The list B_N is linearly independent, so it can be extended to an ordered basis $B_V = (v_1, \dots, v_k, v_{k+1}, \dots, v_n)$ of V (so $\dim(V) = n$). Define $y_i = T(v_{k+i})$ for $1 \leq i \leq n - k$. If we show (y_1, \dots, y_{n-k}) is an ordered basis of R , then $\dim(R) = n - k$ and the statement of the theorem will follow.

Let us first check that the list of y_i 's spans R . Given an arbitrary $z \in R = \text{img}(T)$, we have $z = T(v)$ for some $v \in V$. Expressing v in terms of the basis B_V , there are scalars

$c_1, \dots, c_n \in F$ with $v = \sum_{i=1}^n c_i v_i$. Applying the linear map T to this expression gives $T(v) = \sum_{i=1}^n c_i T(v_i)$. But $T(v_i) = 0$ for $1 \leq i \leq k$ (since these v_i 's are in the kernel of T), and $T(v_i) = y_{i-k}$ for $k < i \leq n$. We therefore get $z = T(v) = \sum_{i=k+1}^n c_i y_{i-k}$, which expresses z as a linear combination of (y_1, \dots, y_{n-k}) .

Next we check that the y_i 's are linearly independent. Assume $d_1, \dots, d_{n-k} \in F$ satisfy $\sum_{j=1}^{n-k} d_j y_j = 0$; we must show every d_j is zero. Since $y_j = T(v_{k+j})$, linearity of T gives $T(\sum_{j=1}^{n-k} d_j v_{k+j}) = 0$. So the vector $u = d_1 v_{k+1} + \dots + d_{n-k} v_n$ is in the kernel of T . As B_N is a basis of N , u can therefore be expressed as some linear combination of v_1, \dots, v_k , say $u = e_1 v_1 + \dots + e_k v_k$ for some $e_i \in F$. Equating the two expressions for u , we get

$$-e_1 v_1 + \dots + (-e_k) v_k + d_1 v_{k+1} + \dots + d_{n-k} v_n = 0.$$

Since (v_1, \dots, v_n) is a linearly independent list, we conclude finally that $-e_1 = \dots = -e_k = d_1 = \dots = d_{n-k} = 0$. This completes the proof of the rank-nullity theorem. The theorem can also be deduced from the fundamental homomorphism theorem for vector spaces (see Exercise 47).

1.9 Summary

Table 1.5 summarizes some of the definitions of algebraic axioms, algebraic systems, subsystems, homomorphisms, and linear algebra concepts discussed in this chapter. Let us also recall the following points.

1. *Examples of Algebraic Structures.* (a) Commutative groups: $\mathbb{Z}, \mathbb{Q}, \mathbb{R}, \mathbb{C}$ under addition; $\mathbb{Z}_n = \{0, 1, \dots, n-1\}$ under addition mod n . (b) Non-commutative groups: $n \times n$ matrices with nonzero determinant under multiplication; bijections on a set X under composition of functions. (c) Commutative rings: $\mathbb{Z}, \mathbb{Q}, \mathbb{R}, \mathbb{C}$ under $+$ and \cdot ; \mathbb{Z}_n (under addition and multiplication mod n); polynomials with coefficients in a field. (d) Non-commutative rings: $n \times n$ matrices with entries in a field. (e) F -vector spaces: F^n with componentwise operations; $m \times n$ matrices with entries in F ; space of F -linear maps between two F -vector spaces. (f) F -algebras: square matrices with entries in F ; space of F -linear operators mapping a given F -vector space to itself; polynomials with coefficients in F .
2. *Product Systems.* If V_1, \dots, V_k are given F -vector spaces, the product set $V = V_1 \times \dots \times V_k = \{(v_1, \dots, v_k) : v_i \in V_i\}$ becomes an F -vector space by defining addition and scalar multiplication componentwise. A similar construction works for other algebraic systems.
3. *Cosets.* If $(G, +)$ is a group with subgroup H , the quotient set G/H is the set of cosets $x + H = \{x + h : h \in H\}$ for $x \in G$. The coset equality theorem states that $x + H = y + H$ iff $(-x) + y \in H$. G is the disjoint union of the distinct cosets of H . When defining functions or operations that act on cosets, one must check that the output is well-defined (independent of the names of the cosets used to compute the answer).
4. *Quotient Systems.* If $(G, +)$ is a commutative group with subgroup H , the set G/H becomes a commutative group with operation $(x + H) + (y + H) = (x + y) + H$ for $x, y \in G$. The identity of G/H is $0_G + H = H$, and the inverse of $x + H$ is $(-x) + H$. A similar construction is possible when G is a group with normal

TABLE 1.5

Summary of Definitions of Algebraic Concepts.

Axiom	Meaning
closure of \star on G	$\forall a, b \in G, a \star b$ lies in G
associativity	$\forall a, b, c, (a \star b) \star c = a \star (b \star c)$
identity e	$\forall a, a \star e = a = e \star a$
inverses	$\forall a, \exists b, a \star b = e = b \star a$
commutativity	$\forall a, b, a \star b = b \star a$
left distributive law	$\forall a, b, c, a \cdot (b + c) = (a \cdot b) + (a \cdot c)$
right distributive law	$\forall a, b, c, (a + b) \cdot c = (a \cdot c) + (b \cdot c)$
Algebraic System	Definition
group (G, \star)	operation \star has closure, associativity, identity, and inverses
commutative group	group where operation is commutative
ring $(R, +, \cdot)$	$(R, +)$ is comm. group; \cdot has closure, associativity, identity, left distributive law, and right distributive law
commutative ring	ring where multiplication is commutative
integral domain	comm. ring with $1 \neq 0$ and no zero divisors
field F	comm. ring where $1 \neq 0$ and every $x \neq 0$ has a mult. inverse
F -vector space $(V, +, s)$	$(V, +)$ is comm. group; $s : F \times V \rightarrow V$ has closure, scalar associativity, two distributive laws, and scalar identity law
R -module	same as vector space, but scalars come from ring R
F -algebra	ring and F -vector space with compatible multiplications
Subsystem	Required Closure Properties
subgroup H of (G, \star)	H closed under \star , identity, and inverses
normal subgroup of G	subgroup closed under conjugation by elements of G
subring S of $(R, +, \cdot)$	S closed under $+, -, \cdot, 0$, and 1
ideal I of $(R, +, \cdot)$	I closed under $+, -, 0$, and left/right mult. by elements of R
subfield K of field F	subring closed under inverses of nonzero elements
subspace of vec. space V	closed under $0, +$, and scalar multiplication
submodule of module	same as subspace, but scalars come from a ring
subalgebra of A	subring and subspace of A
Homomorphism	What Must Be Preserved
group hom.	group operation (hence also identity, inverses, powers)
ring hom.	ring $+$, ring \cdot , multiplicative identity
vec. space hom.	vector addition, scalar multiplication
module hom.	vector addition, scalar multiplication
algebra hom.	ring $+$, ring \cdot , mult. identity, scalar mult.
Lin. Alg. Concept	Definition
spanning list in V	every $v \in V$ is linear combination of vectors on list
linearly independent list	only the zero linear combination of vectors on list gives zero (i.e., $\sum_i c_i v_i = 0$ implies all $c_i = 0$)
ordered basis of V	linearly independent list that spans V
spanning set in V	every $v \in V$ is finite linear combination of vectors from set
linearly independent set	every finite list of distinct elements of set is linearly independent
basis of V	linearly independent set that spans V
dimension of V	size of any basis (or ordered basis) of V

subgroup H ; when G is a ring with ideal H ; when G is a vector space (or module) with subspace (or submodule) H ; and when G is an algebra with subalgebra H . For example, in the F -vector space case, scalar multiplication on G/H is given by $c(x + H) = (cx) + H$ for $c \in F$ and $x \in G$.

5. *Facts about Independence, Spanning, and Bases.* Every F -vector space V has a basis (possibly infinite). Any two bases of V have the same cardinality. Any linearly independent subset of V can be extended to a basis. Any spanning subset of V can be shrunk to a basis. Subsets of linearly independent sets are still linearly independent, and supersets of spanning sets still span. A linearly independent subset of V can never be larger than a spanning set. For a linear map $T : V \rightarrow W$ between finite-dimensional F -vector spaces, $\dim(V) = \dim(\ker(T)) + \dim(\text{img}(T))$ (rank-nullity theorem).
-

1.10 Exercises

1. (a) Prove that a binary operation p on a set S can have at most one identity element. (b) Suppose p is an associative binary operation on a set S with identity element e . Show that each $a \in S$ can have at most one inverse relative to p .
2. Explain why the following sets and binary operations are *not* groups by pointing out a group axiom that fails to hold. (a) S is the set of nonnegative integers, $p(x, y) = x + y$ for $x, y \in S$. (b) $S = \{-1, 0, 1\}$, $p(x, y) = x + y$ for $x, y \in S$. (c) $S = \mathbb{Z}$, $p(x, y) = 4$ for all $x, y \in S$. (d) $S = \mathbb{Z}$, $p(x, y) = (x + y) \bmod 10$ for all $x, y \in S$. (e) $S = \{0, 1, 2, 3, 4\}$, $p(0, x) = x = p(x, 0)$ for all $x \in S$, and $p(x, y) = 0$ for all nonzero $x, y \in S$.
3. Let X be a fixed subset of \mathbb{R} , and let $p(x, y) = \min(x, y)$ for all $x, y \in X$. For each of the axioms for a commutative group, find conditions on X that are necessary and sufficient for (X, p) to satisfy that axiom. For which X is (X, p) a commutative group?
4. (a) For fixed $n \geq 1$, prove that the set of $n \times n$ matrices with real entries is a commutative group under the operation of matrix addition. (b) Is the set of matrices in (a) a group under matrix multiplication? Explain.
5. Decide (with proof) which of the following sets and binary operations are groups.
 - (a) S is the set of injective (one-to-one) functions mapping \mathbb{R} to \mathbb{R} , with $p(f, g) = f \circ g$ (composition of functions) for $f, g \in S$. (b) S is the set of surjective (onto) functions mapping \mathbb{R} to \mathbb{R} , with $p(f, g) = f \circ g$ for $f, g \in S$. (c) S is the set of injective functions mapping \mathbb{Z}_5 to \mathbb{Z}_5 , with $p(f, g) = f \circ g$ for $f, g \in S$. (d) S is the set of all functions from \mathbb{R} to \mathbb{R} , and for $f, g \in S$, $p(f, g)$ is the function sending each $x \in \mathbb{R}$ to $f(x) + g(x)$. (e) S is the set of all functions from \mathbb{R} to \mathbb{R} , and for $f, g \in S$, $p(f, g)$ is the function sending each $x \in \mathbb{R}$ to $f(x)g(x)$.
6. Prove: for all $n \geq 1$, \mathbb{Z}_n is a commutative group under \oplus (addition mod n).
7. (a) Let (G, \star) be a group. Prove the *cancellation laws*: for all $a, x, y \in G$, $a \star x = a \star y$ implies $x = y$, and $x \star a = y \star a$ implies $x = y$. Point out where each of the four group axioms is used in the proof. (b) Give a specific example of a ring $(G, +, \star)$ and nonzero $a, x, y \in G$ for which the implications in part (a) are false. (c) If $(G, +, \star)$ is a field, are the cancellation laws in (a) valid? Prove or give a

- counterexample. (d) Find and prove a version of the multiplicative cancellation law that is valid in an integral domain.
8. (a) *Sudoku Theorem.* Let (G, \star) be a group consisting of the n distinct elements x_1, \dots, x_n . Prove: for all $a \in G$, the list $a \star x_1, \dots, a \star x_n$ is a rearrangement of the list x_1, \dots, x_n ; similarly for $x_1 \star a, \dots, x_n \star a$. (b) *Exponent Theorem.* Let (G, \star) be a finite commutative group of size n . Prove: for all $a \in G$, $a^n = e$ (the identity of G). Do this by letting $G = \{x_1, \dots, x_n\}$ and evaluating the product $(a \star x_1) \star (a \star x_2) \star \dots \star (a \star x_n)$ in two ways. (The theorem holds for non-commutative groups too, but the proof is harder.)
 9. Let A and B be normal subgroups of a group G . (a) Prove: if $A \cap B = \{e_G\}$, then $ab = ba$ for all $a \in A$ and $b \in B$. (Study $aba^{-1}b^{-1}$.) (b) Prove $AB = \{ab : a \in A, b \in B\}$ is a subgroup of G . Is AB normal in G ?
 10. Fix $n \geq 1$. Let \oplus and \otimes denote addition mod n and multiplication mod n , respectively. (a) Prove: for all $a, b \in \mathbb{Z}_n$, $a \oplus b \in \mathbb{Z}_n$ and $a \oplus b = a + b - qn$ for some $q \in \mathbb{Z}$. (b) Prove: for all $a, b \in \mathbb{Z}_n$, $a \otimes b \in \mathbb{Z}_n$ and $a \otimes b = ab - sn$ for some $s \in \mathbb{Z}$. (c) Prove: for all $c, d \in \mathbb{Z}_n$, $c = d$ iff $c - d = tn$ for some $t \in \mathbb{Z}$. (d) Prove $(\mathbb{Z}_n, \oplus, \otimes)$ is a commutative ring. [Hint: Use (a) and (b) to eliminate \oplus and \otimes from each side of each ring axiom; then use (c) to see that the two sides must be equal.]
 11. For any ring $(R, +, \cdot)$, let R^* (the set of *units* of R) be the set of all $x \in R$ for which there exists $y \in R$ with $x \cdot y = 1_R = y \cdot x$. (a) Prove that (R^*, \cdot) is a group, and show R^* is commutative if R is commutative. (b) Describe R^* for each of these rings: $R = \mathbb{Z}$; $R = \mathbb{R}$; R is any field; $R = \mathbb{Z}_{12}$; $R = M_3(\mathbb{R})$; $R = \mathbb{R}[x]$.
 12. Let S be the set of all functions $f : \mathbb{R} \rightarrow \mathbb{R}$. For $f, g \in S$, let $f + g$ be the function sending each $x \in \mathbb{R}$ to $f(x) + g(x)$, and let $f \cdot g$ be the function sending each $x \in \mathbb{R}$ to $f(x) \cdot g(x)$. (a) Prove $(S, +, \cdot)$ is a commutative ring. Is this a field? Is this an integral domain? (b) Consider $(S, +, \circ)$ where \circ is composition of functions. Which axioms for a commutative ring hold?
 13. (a) Prove that every field is an integral domain. (b) Prove that every *finite* integral domain R is a field. [Hint: For nonzero $a \in R$, study whether the map $L_a : R \rightarrow R$, given by $L_a(x) = a \cdot x$ for $x \in R$, is one-to-one or onto.] (c) Prove: for all $n \geq 1$, \mathbb{Z}_n is a field iff \mathbb{Z}_n is an integral domain iff n is prime.
 14. Let F be a field. (a) Fix $k \geq 1$. Carefully verify that F^k (with componentwise operations) is a vector space over F . (b) More generally, verify the vector space axioms for the direct product $V_1 \times V_2 \times \dots \times V_n$ of F -vector spaces V_1, V_2, \dots, V_n .
 15. Let V be the set of integers under ordinary addition. For $c \in \mathbb{R}$ and $v \in V$, let $s(c, v) = cv$ be the ordinary product of the real numbers c and v . Which axioms for a real vector space hold for $(V, +, s)$?
 16. For $V = \mathbb{R}$, define scalar multiplication by $s(c, v) = 0$ for all $c, v \in \mathbb{R}$. Show that $(V, +, s)$ satisfies all the axioms for a real vector space except the identity axiom for scalar multiplication.
 17. For $V = \mathbb{R}$, define scalar multiplication by setting $s(c, v) = v$ for all $c, v \in \mathbb{R}$. Show that $(V, +, s)$ satisfies all the axioms for a real vector space except for distributivity of scalar addition.
 18. For $V = \mathbb{R}$ and $F = \mathbb{C}$, define scalar multiplication $s : F \times V \rightarrow V$ by setting $s(a + ib, v) = av$ for all $a, b, v \in \mathbb{R}$. Show that $(V, +, s)$ satisfies all the axioms for a complex vector space except for associativity of scalar multiplication.

19. (a) Prove that for the field $F = \mathbb{Z}_p$ (where p is prime), distributivity of vector addition follows from the other axioms for an F -vector space. (b) Prove that the result in (a) also holds for the field $F = \mathbb{Q}$. (c) Can you find an example of a field F and a structure $(V, +, s)$ satisfying all the F -vector space axioms except distributivity of vector addition?
20. For $V = \mathbb{R}$, define scalar multiplication s by setting $s(c, v) = c^2 v$ for all $c, v \in \mathbb{R}$. Which axioms for a real vector space hold for $(V, +, s)$?
21. For $V = \mathbb{R}$, define scalar multiplication s by setting $s(c, v) = c$ for all $c, v \in \mathbb{R}$. Which axioms for a real vector space hold for $(V, +, s)$?
22. Let $V = \mathbb{R}^+$, the set of positive real numbers. Define $p(v, w) = vw$ for $v, w \in V$, and define $s(c, v) = v^c$ for $c \in \mathbb{R}$ and $v \in V$. Which axioms for a real vector space hold for (V, p, s) ?
23. Check in detail that the examples of subsystems in the last paragraph of §1.4 really do have the required closure properties.
24. Let W be the set of matrices of the form $\begin{bmatrix} a & 0 \\ b & 0 \end{bmatrix}$ for some $a, b \in \mathbb{R}$. Is W a subgroup of $M_2(\mathbb{R})$? a subring? a left ideal? a right ideal? an ideal? a subspace? a subalgebra? Explain.
25. Give an example of a ring R and a right ideal I in R that is not a two-sided ideal of R .
26. If possible, give an example of a subset of \mathbb{Z} that is closed under addition and inverses, yet is not a subgroup of $(\mathbb{Z}, +)$.
27. For any field F and $k \geq 1$, show that $W = \{(t, t, \dots, t) : t \in F\}$ is a subspace of the F -vector space F^k . Is an analogous result true for groups or for rings?
28. Let G be the group of invertible 2×2 matrices with real entries. (a) Give three different examples of normal subgroups of G . (b) Give an example of a non-normal subgroup of G .
29. Prove that $\mathbb{Z}[i] = \{a + bi : a, b \in \mathbb{Z}\}$ is a subring of \mathbb{C} . Is this a subfield of \mathbb{C} ?
30. Given rings R_1, \dots, R_n , carefully prove that the product set $R = R_1 \times \dots \times R_n$ is a ring (with componentwise operations). In particular, what are 0_R and 1_R ?
31. Prove that the direct product of two or more fields is never a field.
32. Suppose G_1 and G_2 are groups, H_1 is a subgroup of G_1 , and H_2 is a subgroup of G_2 . (a) Prove $H_1 \times H_2$ is a subgroup of $G_1 \times G_2$, which is normal if H_1 is normal in G_1 and H_2 is normal in G_2 . (b) State and prove analogous results for subrings and ideals in a product ring and for subspaces in a product vector space. (c) Give an example of vector spaces V_1 and V_2 and a subspace W of $V_1 \times V_2$ that is not of the form $W_1 \times W_2$ for any choice of W_1, W_2 .
33. Let (G, \star) be a group with subgroup H . (a) *Left coset equality theorem.* Prove: for all $a, b \in G$, the following conditions are all equivalent: $a \star H = b \star H$; $a \in b \star H$; $b \in a \star H$; $a = b \star h$ for some $h \in H$; $b = a \star k$ for some $k \in H$; $b^{-1} \star a \in H$; $a^{-1} \star b \in H$. (b) Formulate and prove an analogous result for right cosets. (c) Prove: H is a normal subgroup of G iff for all $a \in G$, $a \star H = H \star a$.
34. Let H be a normal subgroup of a group (G, \star) . (a) Verify that the binary operation $p : G/H \times G/H \rightarrow G/H$, given by $p(a \star H, b \star H) = (a \star b) \star H$ for all $a, b \in G$, is well-defined. Indicate where your proof uses normality of H . (b) Verify that $(G/H, p)$ is a group.

35. Let I be an ideal in a ring $(R, +, \cdot)$. (a) Prove that the binary operation $(a + I) \bullet (b + I) = (a \cdot b) + I$ for $a, b \in R$ is well-defined. (b) Prove that $(R/I, +, \bullet)$ is a ring, which is commutative if R is commutative.
36. Let $W = \{(t, -t) : t \in \mathbb{R}\}$, which is a subspace of the real vector space \mathbb{R}^2 . Draw a picture of the set \mathbb{R}^2/W . Calculate $[(1, 2) + W] + [(0, -1) + W]$, illustrate this calculation on your picture, and give three different names for the answer.
37. Suppose G and K are groups, and $T : G \rightarrow K$ is a group homomorphism. Prove: for all $x \in G$ and all $n \in \mathbb{Z}$, $T(x^n) = T(x)^n$.
38. Let X , Y , and Z be vector spaces over a field F , and let $T : X \rightarrow Y$ and $U : Y \rightarrow Z$ be F -linear maps. (a) Prove $U \circ T : X \rightarrow Z$ is F -linear. (b) Prove: if T is a vector space isomorphism, so is T^{-1} . (c) Prove $\ker(T)$ is a subspace of X , and $\text{img}(T)$ is a subspace of Y . (d) Prove: $\ker(T) \subseteq \ker(U \circ T)$. Find a condition under which equality holds. (e) Prove: $\text{img}(U \circ T) \subseteq \text{img}(U)$. Can you find conditions under which equality will hold?
39. State and prove the fundamental homomorphism theorem for vector spaces over a field F . [Hint: Your proof will be quite short if you invoke the fundamental homomorphism theorem for additive groups.]
40. Given a field F and an integer $k > 0$, define vectors $e_1, \dots, e_k \in F^k$ by letting e_i have 1_F in position i and 0_F in all other positions. (a) Prove that (e_1, e_2, \dots, e_k) is an ordered basis of F^k . (b) Define $f_i = e_1 + e_2 + \dots + e_i$ for $1 \leq i \leq k$. Prove that (f_1, f_2, \dots, f_k) is an ordered basis of F^k .
41. For each real vector space of matrices, find a basis for the space and compute its dimension. (a) the set of all real $n \times n$ matrices; (b) the set of all complex $n \times n$ matrices; (c) the set of all real upper-triangular $n \times n$ matrices; (d) the set of all real symmetric $n \times n$ matrices.
42. Let $S \subseteq T \subseteq V$, where V is an F -vector space. Use the definitions to prove the following results. (a) If S spans V , then T spans V . (b) If T is F -linearly independent, then S is F -linearly independent. (c) V spans V . (d) \emptyset is linearly independent.
43. Let V and W be F -vector spaces, and let $T : V \rightarrow W$ be an F -linear map. (a) Prove: if U is a subspace of V , then $\dim(U) \leq \dim(V)$. (b) Prove: if U is a subspace of V and $\dim(U) = \dim(V) < \infty$, then $U = V$. Must this conclusion hold if V is infinite-dimensional? (c) Prove: $\dim(\text{img}(T)) \leq \min(\dim(V), \dim(W))$. (d) Assuming V and W are finite-dimensional, find and prove inequalities relating $\dim(\ker(T))$ to $\dim(V)$ and $\dim(W)$.
44. *Subspace Generated by Vectors.* Let V be an F -vector space, and let S be a list or set of vectors in V . Prove that the set W of all finite F -linear combinations of vectors in S is a subspace of V . (This subspace is called the *subspace of V generated by S* .) Prove also that for every subspace Z of V that contains S , $W \subseteq Z$.
45. Let V and W be F -vector spaces, let $T : V \rightarrow W$ be an F -linear map, and let U be a subspace of V . (a) Prove that $T[U] = \{T(u) : u \in U\}$ is a subspace of W . (b) Prove: if a list (u_1, \dots, u_k) spans U , then the list $(T(u_1), \dots, T(u_k))$ spans $T[U]$.
46. Let V and W be F -vector spaces, and let $T : V \rightarrow W$ be an F -linear map. (a) Prove: T is one-to-one iff for every linearly independent list (v_1, \dots, v_k) in V , the list $(T(v_1), \dots, T(v_k))$ is linearly independent in W . (b) Prove: T is onto iff

- for every spanning set S of V , $T[S] = \{T(x) : x \in S\}$ spans W . (c) Prove: T is an isomorphism iff for every basis B of V , $T[B]$ is a basis of W .
47. (a) Let V be a finite-dimensional F -vector space with subspace W . Prove that $\dim(V/W) = \dim(V) - \dim(W)$. [Hint: Extend a basis of W to a basis of V , and prove that the cosets of the new basis vectors form a basis of V/W .] (b) Use (a) and the fundamental homomorphism theorem for vector spaces (Exercise 39) to reprove the rank-nullity theorem.
48. *Endomorphism Ring of a Commutative Group.* Given a commutative group $(M, +)$, let $\text{End}(M)$ be the set of all group homomorphisms $f : M \rightarrow M$. $\text{End}(M)$ is called the *endomorphism ring of M* . Define the *sum* of $f, g \in \text{End}(M)$ to be the function $f + g : M \rightarrow M$ given by $(f + g)(x) = f(x) + g(x)$ for all $x \in M$. Define the *product* of $f, g \in \text{End}(M)$ to be the composition $f \circ g$, given by $(f \circ g)(x) = f(g(x))$ for all $x \in M$. (a) Verify the ring axioms for $\text{End}(M)$ with these operations. (b) Let $\text{Fun}(M)$ be the set of all functions $f : M \rightarrow M$. Using the sum and product defined above, determine which ring axioms hold in $\text{Fun}(M)$.

This page intentionally left blank

Permutations

This chapter gives the basic definitions and facts about permutations that are needed to discuss determinants and multilinear algebra. First we discuss how composition of functions confers a group structure on the set S_n of permutations of $\{1, 2, \dots, n\}$. Next we introduce a way of visualizing a function using a directed graph, which leads to a description of permutations in terms of disjoint directed cycles. We use this description to obtain some algebraic factorizations of permutations in the group S_n . The chapter concludes by studying inversions of functions and permutations, which give information about how many steps it takes to sort a list into increasing order. We use inversions to define the sign of a permutation, which will play a critical role in our subsequent treatment of determinants (Chapter 5).

2.1 Symmetric Groups

This section assumes familiarity with the definition of a group (§1.1). For each positive integer n , let $[n]$ denote the finite set $\{1, 2, \dots, n\}$. The *symmetric group* S_n is defined to be the set of all bijective functions $f : [n] \rightarrow [n]$, with composition of functions as the binary operation. Recall that a function $f : X \rightarrow Y$ is a bijection iff f is *one-to-one* (for all $x, z \in X$, $f(x) = f(z)$ implies $x = z$) and f is *onto* (for each $y \in Y$ there is $x \in X$ with $f(x) = y$). For a function f mapping the *finite* set $[n]$ into itself, one may check that f is one-to-one iff f is onto.

Recall that the *composition* of two functions $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ is defined to be the function $g \circ f : X \rightarrow Z$ given by $(g \circ f)(x) = g(f(x))$ for $x \in X$. The composition of bijections is a bijection, so S_n is closed under the composition operation. Composition of functions is always associative, so associativity holds. The *identity function* on $[n]$, given by $\text{id}(x) = x$ for $x \in [n]$, is a bijection and hence belongs to S_n . Since $f \circ \text{id} = f = \text{id} \circ f$ for all $f \in S_n$, S_n has an identity element relative to \circ . Finally, each $f \in S_n$ has an inverse function $f^{-1} : [n] \rightarrow [n]$, which is also a bijection and satisfies $f \circ f^{-1} = \text{id} = f^{-1} \circ f$. Therefore, (S_n, \circ) is a group. S_n is not a commutative group except when $n \leq 2$.

Elements of S_n are called *permutations* of n objects. To explain this terminology, note that any function $f : [n] \rightarrow [n]$ is completely determined by the list $[f(1), f(2), \dots, f(n)]$. We refer to this list as the *one-line form* for f . If f belongs to S_n , so that $i \neq j$ implies $f(i) \neq f(j)$, then this list is a rearrangement (or “permutation”) of the list $[1, 2, \dots, n]$.

For example, the function $f : [4] \rightarrow [4]$ given by $f(1) = 3$, $f(2) = 2$, $f(3) = 4$, and $f(4) = 1$ has one-line form $f = [3, 2, 4, 1]$. Given the one-line form $g = [4, 1, 3, 2] \in S_4$, we must have $g(1) = 4$, $g(2) = 1$, $g(3) = 3$, and $g(4) = 2$. Then $f \circ g = [1, 3, 4, 2]$ since $f \circ g(1) = f(g(1)) = f(4) = 1$, $f(g(2)) = f(1) = 3$, etc. On the other hand, $g \circ f = [3, 1, 2, 4]$.

Using one-line forms, we can show that the group S_n has exactly $n!$ elements. For, there are n ways to choose $f(1)$; then $n - 1$ ways to choose $f(2) \in [n] \sim \{f(1)\}$; then $n - 2$ ways to choose $f(3) \in [n] \sim \{f(1), f(2)\}$; and so on. Finally, there is $n - (n - 1) = 1$ way to choose $f(n)$, so that there are $n \times (n - 1) \times \dots \times 1 = n!$ permutations in S_n .

2.2 Representing Functions as Directed Graphs

Let $f : [n] \rightarrow [n]$ be any function. We can create a graphical representation of f by drawing n dots labeled 1 through n , and then drawing an arrow from i to $f(i)$ for all $i \in [n]$. If $i = f(i)$, this arrow is a loop pointing from i to itself. See Figure 2.1 for an example where $n = 23$ and

$$f = [10, 22, 7, 11, 15, 19, 19, 12, 22, 12, 11, 1, 6, 11, 5, 8, 22, 9, 21, 11, 3, 3, 2].$$

Note that we can recover the function f from its directed graph.

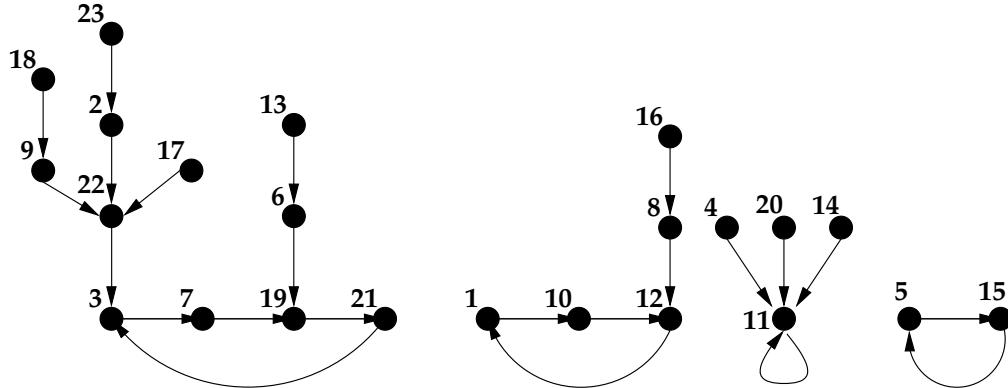


FIGURE 2.1
Directed Graph for a Function.

Each node in the graph of f has exactly one arrow emanating from it, since f is a function. Thus, starting from some given node x_0 , there is a well-defined path through the graph obtained by repeatedly following the arrows for f . This path visits the nodes $x_0, x_1 = f(x_0), x_2 = f(f(x_0)), x_3 = f(f(f(x_0))),$ and so on. Since there are only finitely many nodes in the graph, we must have $x_i = x_j$ for some $i < j < \infty$. If i and j are the smallest indices for which this is true, we see that the path we followed ends in a directed “cycle” (consisting of the nodes $x_i, x_{i+1}, \dots, x_j = x_i$) and begins in a “tail” x_0, \dots, x_{i-1} that feeds into the cycle at position x_i . If $i = 0$, there is no tail, and the path we are following lies completely on the cycle. If $j = i + 1$, the cycle consists of a single loop edge based at the node x_i .

2.3 Cycle Decompositions of Permutations

Suppose that $f : [n] \rightarrow [n]$ is a permutation (i.e., a bijection). What does the directed graph for f look like in this case? Since f is onto, there is at least one arrow entering each node in the graph. Since f is one-to-one, there is at most one arrow entering each node in the graph. Therefore, for the path described at the end of the last section, there can be no tail feeding into the cycle at position x_i — otherwise, the arrows starting at x_{i-1} and x_{j-1} both land at x_i , which is impossible. Hence, every path obtained by following arrows through

the graph of f is a cycle, possibly involving just one node. Every node belongs to such a cycle, since every node has an arrow leaving it. Different cycles visit disjoint sets of nodes, since otherwise there would be two arrows entering the same node. So, the directed graph representing a permutation consists of a disjoint union of cycles.

It follows that we can define a permutation f of $[n]$ by listing the elements in each of the cycles of f 's directed graph, in the order they appear along the cycle. This is called a *cycle decomposition* of f . The notation (i_1, i_2, \dots, i_k) represents a cycle that visits the k distinct nodes i_1, i_2, \dots, i_k in this order and then returns to i_1 . Observe that (i_2, \dots, i_k, i_1) and $(i_3, \dots, i_k, i_1, i_2)$ are alternate notations for the cycle just mentioned, but (i_k, \dots, i_2, i_1) is a different cycle (for $k \geq 3$) since the direction of the edges matters. If the graph of f has multiple cycles, we juxtapose the notation for each individual cycle to get a cycle decomposition of f . For example, let $n = 6$ and $f = [4, 2, 5, 6, 3, 1]$ in one-line form. The directed graph for f is shown in Figure 2.2, and a cycle decomposition of f is $f = (1, 4, 6)(2)(3, 5)$. We can obtain other cycle decompositions of f by listing the cycles in the directed graph of f in different orders, or by beginning the traversal of each cycle in different places.

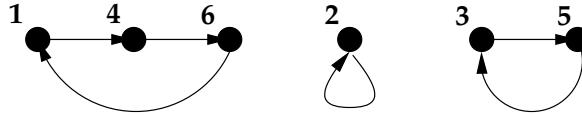


FIGURE 2.2

Directed Graph of a Permutation.

A cycle containing k nodes is called a *k -cycle*. By convention, we are permitted to omit any or all 1-cycles from the cycle decomposition of a permutation. For example, f could also be written $f = (1, 4, 6)(3, 5)$ or $(5, 3)(6, 1, 4)(2)$. For more examples, note that individual cycles of any length now define permutations of $[n]$ by restoring the omitted 1-cycles. For instance, we have permutations $f_1 = (1, 4, 6)$, $f_2 = (2) = \text{id}$, and $f_3 = (3, 5)$, which are given in one-line form by $f_1 = [4, 2, 3, 6, 5, 1]$, $f_2 = [1, 2, 3, 4, 5, 6]$, and $f_3 = [1, 2, 5, 4, 3, 6]$.

We took care to use *square brackets* in the definition of one-line form to distinguish this notation from the notation for an n -cycle. For example, $f = [4, 2, 5, 6, 3, 1]$ is different from the 6-cycle $(4, 2, 5, 6, 3, 1) = (1, 4, 2, 5, 6, 3)$. The one-line form of this 6-cycle is $[4, 5, 1, 2, 6, 3] \neq f$.

We can compose two permutations using their cycle decompositions. Given $f = (1, 4, 6)(2)(3, 5)$ and $g = (2, 6, 4, 3)(1, 5)$, let us find $f \circ g$ and $g \circ f$. For $f \circ g$, we could compute $f(g(1)) = f(5) = 3$, $f(g(2)) = f(6) = 1$, $f(g(3)) = f(2) = 2$, etc., leading to the one-line form $f \circ g = [3, 1, 2, 5, 4, 6]$. Converting this to a cycle decomposition, we have $f \circ g = (1, 3, 2)(4, 5)(6)$. On the other hand, let us directly compute a cycle decomposition of $g \circ f$ without finding the one-line form first. Node 1 belongs to some cycle of $g \circ f$; indeed, $g(f(1)) = g(4) = 3$, so 1 is followed on its cycle by 3. Next, $g(f(3)) = g(5) = 1$, so 3 is followed on the cycle by 1, and we see that the 2-cycle $(1, 3)$ is part of the directed graph of $g \circ f$. Next we consider node 2, and note $g(f(2)) = g(2) = 6$, then $g(f(6)) = g(1) = 5$, then $g(f(5)) = g(3) = 2$, so we have the cycle $(2, 6, 5)$. Finally, the cycle involving node 4 is found by noting $g(f(4)) = g(6) = 4$. Thus, $g \circ f = (1, 3)(2, 6, 5)(4)$, which could also be written $(6, 5, 2)(3, 1)$, or in many other ways. If needed, we can now find the one-line form $g \circ f = [3, 6, 1, 4, 2, 5]$.

Let $f = (1, 4, 6)(2)(3, 5)$, $f_1 = (1, 4, 6)$, $f_2 = (2)$, and $f_3 = (3, 5)$ as above. By a

computation similar to the one given in the last paragraph, one sees that

$$f = (1, 4, 6)(2)(3, 5) = (1, 4, 6) \circ (2) \circ (3, 5) = f_1 \circ f_2 \circ f_3.$$

Intuitively, the action of f on $[n]$ can be accomplished by first applying f_3 , which moves only the elements on the cycle $(3, 5)$; then applying f_2 , which happens to be the identity map; and then applying f_1 , which moves only the elements on the cycle $(1, 4, 6)$. Similarly, $f = f_1 \circ f_3 \circ f_2 = f_2 \circ f_1 \circ f_3 = f_3 \circ f_1 \circ f_2$, etc. Intuitively, we achieve the same net effect no matter what order we apply the functions f_1 , f_2 , and f_3 , which translates into the algebraic statement that these three elements of S_n commute with each other.

2.4 Composition of Cycles

We now generalize the remarks at the end of the last section. Consider two cycles $g = (a_1, \dots, a_k)$ and $h = (b_1, \dots, b_m)$, which we view as elements of S_n . We say g and h are *disjoint* cycles iff $\{a_1, \dots, a_k\}$ and $\{b_1, \dots, b_m\}$ are disjoint sets. If g and h are disjoint, we claim that these two cycles commute in the group S_n , i.e., $g \circ h = h \circ g$. To verify this claim carefully, first recall that two functions $p, q : [n] \rightarrow [n]$ are *equal* iff $p(x) = q(x)$ for all $x \in [n]$. In the situation at hand, we have $p = g \circ h$ and $q = h \circ g$. Take an arbitrary $x \in [n]$, and consider three cases. First, if $x = a_i$ for some i , then

$$g \circ h(x) = g(h(a_i)) = g(a_i) = h(g(a_i)) = h \circ g(x)$$

since $h(a_j) = a_j$ for all j , and $g(a_i)$ equals some a_j . Second, if $x = b_i$ for some i , then

$$g \circ h(x) = g(h(b_i)) = h(b_i) = h(g(b_i)) = h \circ g(x)$$

since $g(b_j) = b_j$ for all j , and $g(b_i)$ equals some b_j . Third, if x is different from all a_i 's and all b_i 's, then

$$g \circ h(x) = x = h \circ g(x).$$

The above discussion tacitly assumed that neither g nor h was the identity permutation id . But, it is immediate that id commutes with every element of S_n .

We now verify a few more identities involving composition of cycles. Suppose $f = C_1 C_2 \cdots C_k$ is any cycle decomposition of $f \in S_n$. Let us check carefully that $f = C_1 \circ C_2 \circ \cdots \circ C_k$. The notation C_i refers both to a cycle in the directed graph of f , and the function in S_n corresponding to that cycle. We must show that the function $C_1 \circ C_2 \circ \cdots \circ C_k$ applied to any $x \in [n]$ yields $f(x)$. If x does not appear in any of the cycles C_i , then all of the permutations C_k, C_{k-1}, \dots, C_1 leave x fixed.¹ Also, f itself must fix x since we are only allowed to omit 1-cycles in a cycle decomposition. Thus, $f(x) = x = C_1 \circ C_2 \circ \cdots \circ C_k(x)$. The other possibility is that x appears in a cycle C_i for some (necessarily unique) i . If $f(x) = y$, then $C_i(x) = y$ because of the way we constructed the cycles from the directed graph of f . All the other cycles C_j do not overlap with C_i , and therefore the functions associated with these other cycles must fix x and y . Starting at x and applying the functions C_k through C_1 in this order, x is always fixed until C_i sends x to y , and then y is always fixed. Hence, $C_1 \circ \cdots \circ C_k(x) = y = f(x)$, as needed.

Next, consider a general k -cycle (a_1, \dots, a_k) , where $k \geq 3$. We show

$$(a_1, \dots, a_k) = (a_1, a_k) \circ (a_1, a_{k-1}) \circ \cdots \circ (a_1, a_3) \circ (a_1, a_2). \quad (2.1)$$

¹We say a function g *fixes* an element x iff $g(x) = x$.

Both sides fix every x unequal to all a_j 's. Both sides send a_1 to a_2 and a_k to a_1 . (Remember that the maps on the right side are applied from right to left.) For $1 < j < k$, the left side sends a_j to a_{j+1} . Reading the other side from right to left, the first $(j-2)$ 2-cycles leave a_j fixed. The next 2-cycle, namely (a_1, a_j) , sends a_j to a_1 . The next 2-cycle, namely (a_1, a_{j+1}) , sends a_1 to a_{j+1} . The remaining 2-cycles fix a_{j+1} , so the net effect is that a_j goes to a_{j+1} .

By entirely analogous arguments, one can rigorously verify the identity

$$(j, i) = (i, j) = (i, i+1) \circ (i+1, i+2) \circ \cdots \circ (j-2, j-1) \circ (j-1, j) \\ \circ (j-2, j-1) \circ \cdots \circ (i+1, i+2) \circ (i, i+1), \quad (2.2)$$

which is valid for all $i, j \in [n]$ with $i < j$.

2.5 Factorizations of Permutations

We now consider various ways of factoring permutations in the group S_n into products of simpler permutations. Here, “product” refers to the group operation in S_n , which is composition of functions. These results are analogous to the fundamental theorem of arithmetic, which states that any positive integer can be written as a product of prime integers.

Factorization into disjoint cycles. Given any permutation $f \in S_n$, our analysis of the directed graph of f shows that f has a cycle decomposition $f = C_1 C_2 \cdots C_k$. We proved above that this cycle decomposition leads to the factorization $f = C_1 \circ C_2 \circ \cdots \circ C_k$, which expresses f as a product of pairwise disjoint cycles in S_n .

To what extent is this factorization *unique*? On one hand, the cycles C_1, \dots, C_k are pairwise disjoint, so they all commute with one another. On the other hand, we are free to add or remove 1-cycles from this factorization, since all 1-cycles represent the identity element of S_n . Disregarding the order of factors and the possible omission of 1-cycles, it can be shown that the factorization into disjoint cycles is unique. Intuitively, this follows because the non-identity cycles in any factorization of f into a product of disjoint cycles must correspond (in some order) to the cycles in the directed graph of f , and these cycles are uniquely determined by f .

Factorizations into transpositions. A *transposition* is a 2-cycle (i, j) . The formula (2.1) shows that any k -cycle with $k \geq 3$ can be factored in S_n as a product of transpositions. Applying this formula to each cycle in a cycle decomposition of $f \in S_n$, we see that *every element of S_n can be written as a product of zero or more transpositions*. Here, an empty product represents the identity permutation, which can also be factored as $\text{id} = (1, 2) \circ (1, 2)$ for $n \geq 2$. Note that the transpositions used in such a factorization need not be disjoint. Furthermore, the factorization of a permutation into transpositions is not unique. For example, given any factorization of any $g \in S_n$ into transpositions, we can always create new factorizations of g by appending pairs of factors of the form $(i, j) \circ (i, j)$. For a less silly example of non-uniqueness, one may check that

$$[4, 1, 2, 3] = (1, 3) \circ (1, 4) \circ (2, 3) = (3, 4) \circ (2, 3) \circ (1, 2).$$

Factorizations into basic transpositions. A *basic transposition* is a transposition of the form $(i, i+1)$ for some $i < n$. The formula (2.2) shows that any transposition can be factored in S_n into a product of basic transpositions. Applying this formula to each transposition in a

factorization of $f \in S_n$ into transpositions, we see that *every element of S_n can be written as a product of zero or more basic transpositions*. These factorizations are not unique; for instance, $(1, 2) \circ (1, 3) = (1, 3, 2) = (2, 3) \circ (1, 2)$. We will see in the next section another proof (and algorithm) showing how $f \in S_n$ can be written as a product of basic transpositions.

2.6 Inversions and Sorting

Given positive integers m and n , consider a function $f : [n] \rightarrow [m]$. An *inversion* of f is a pair (i, j) such that $1 \leq i < j \leq n$ and $f(i) > f(j)$. In terms of the one-line form $[f(1), f(2), \dots, f(n)]$, an inversion of f is a pair of positions (i, j) , not necessarily adjacent, such that $f(i)$ and $f(j)$ appear out of order. We let $\text{Inv}(f)$ denote the set of inversions of f , and we let $\text{inv}(f) = |\text{Inv}(f)|$ be the number of inversions of f . These concepts apply, in particular, to permutations $f \in S_n$. For example, given $f = [3, 1, 4, 2]$, one inversion of f is $(1, 2)$, since $f(1) = 3 > 1 = f(2)$. Listing all the inversions, we find that $\text{Inv}(f) = \{(1, 2), (1, 4), (3, 4)\}$ and $\text{inv}(f) = 3$. For another example, $g = [1, 1, 2, 1, 1, 2]$ has $\text{Inv}(g) = \{(3, 4), (3, 5)\}$ and $\text{inv}(g) = 2$.

Given $f : [n] \rightarrow [m]$, suppose we want to sort the list $[f(1), f(2), \dots, f(n)]$ into weakly increasing order by making a series of “basic transposition moves.” By definition, a *basic transposition move* consists of switching two adjacent list elements that are strictly out of order; in other words, we are permitted to switch $f(i)$ and $f(i+1)$ iff $f(i) > f(i+1)$. How many such moves does it take to sort the list? We claim that, regardless of which moves are made in what order, the list will always be sorted in exactly $\text{inv}(f)$ moves. This claim follows from the following three assertions: (1) The list is fully sorted iff $\text{inv}(f) = 0$. (2) If the list is not fully sorted, there exists at least one permissible basic transposition move. (3) If g is obtained from f by applying any one basic transposition move, then $\text{inv}(g) = \text{inv}(f) - 1$. Assertions (1) and (2) are readily verified. To check (3), suppose $f(i) > f(i+1)$, and we get g from f by switching these two elements in the one-line form. By definition of inversions, we see that the set $\text{Inv}(f)$ can be transformed into $\text{Inv}(g)$ by deleting the inversion $(i, i+1)$, and then replacing i by $i+1$ and $i+1$ by i in all other inversion pairs in which i or $i+1$ appears. We see from this that $\text{Inv}(g)$ has one less element than $\text{Inv}(f)$, so that (3) holds. For example, if $f = [3, 1, 4, 2]$ is transformed into $g = [3, 1, 2, 4]$ by switching the last two elements, then $\text{Inv}(f) = \{(1, 2), (1, 4), (3, 4)\}$ is transformed into $\text{Inv}(g) = \{(1, 2), (1, 3)\}$ by the process just described.

This sorting process provides a nice algorithm for factoring permutations $f \in S_n$ into products of basic transpositions. Suppose g is obtained from f by switching $f(i)$ and $f(i+1)$ in the one-line form of f . One may check that $g = f \circ (i, i+1)$, so that (in the case $f(i) > f(i+1)$) the basic transposition move switching the entries in positions i and $i+1$ can be accomplished algebraically in S_n by multiplying f on the right by $(i, i+1)$. Sorting f to $\text{id} = [1, 2, \dots, n]$ by a sequence of such moves will therefore give a formula of the form

$$f \circ (i_1, i_1 + 1) \circ (i_2, i_2 + 1) \circ \cdots \circ (i_k, i_k + 1) = \text{id},$$

where $k = \text{inv}(f)$. Solving for f , we find that $f = (i_k, i_k + 1) \circ \cdots \circ (i_2, i_2 + 1) \circ (i_1, i_1 + 1)$. In our running example where $f = [3, 1, 4, 2]$, we can write $f \circ (3, 4) \circ (1, 2) \circ (2, 3) = \text{id}$, and hence $f = (2, 3) \circ (1, 2) \circ (3, 4)$. This argument proves that *any $f \in S_n$ can be factored into the product of $\text{inv}(f)$ basic transpositions*.

2.7 Signs of Permutations

For any function $f : [n] \rightarrow [n]$, define the *sign* of f to be $\text{sgn}(f) = (-1)^{\text{inv}(f)}$ if f is a bijection, and $\text{sgn}(f) = 0$ otherwise. For example, if $f = \text{id} = [1, 2, \dots, n]$, then $\text{inv}(f) = 0$ and $\text{sgn}(f) = 1$. If h is the basic transposition $(i, i + 1)$, then

$$h = [1, \dots, i - 1, i + 1, i, i + 2, \dots, n],$$

so that $\text{Inv}(h) = \{(i, i + 1)\}$, $\text{inv}(h) = 1$, and $\text{sgn}(h) = -1$.

A fundamental theorem about sgn is that it is a group homomorphism from S_n into the multiplicative group $\{-1, 1\}$, i.e.,

$$\text{sgn}(f \circ h) = \text{sgn}(f) \cdot \text{sgn}(h) \text{ for all } f, h \in S_n. \quad (2.3)$$

First, observe that this relation holds when $h = \text{id}$. Next, consider the special case where h is a basic transposition $(i, i + 1)$. Fix $f \in S_n$ and set $g = f \circ h = f \circ (i, i + 1)$. We have already observed that the one-line form for g is obtained from the one-line form for f by switching the entries $f(i)$ and $f(i + 1)$ in positions i and $i + 1$. Hence, if $f(i) > f(i + 1)$, then g is obtained from f by applying a basic transposition move. By assertion (3) in §2.6, $\text{inv}(g) = \text{inv}(f) - 1$, and so

$$\text{sgn}(g) = -\text{sgn}(f) = \text{sgn}(f) \cdot \text{sgn}((i, i + 1)) = \text{sgn}(f) \cdot \text{sgn}(h).$$

In the case where $f(i) < f(i + 1)$, we see that f is obtained from g by applying a basic transposition move. Applying assertion (3) with f and g interchanged, we see that $\text{inv}(g) = \text{inv}(f) + 1$, and again

$$\text{sgn}(g) = -\text{sgn}(f) = \text{sgn}(f) \cdot \text{sgn}((i, i + 1)) = \text{sgn}(f) \cdot \text{sgn}(h).$$

Therefore, (2.3) holds when h is any basic transposition.

To prove (2.3) for arbitrary $h \in S_n$, we can write $h = h_1 \circ \dots \circ h_s$ where each h_j is a basic transposition. We now argue by induction on s . If $s \leq 1$, we have shown that (2.3) holds for all $f \in S_n$. Assuming $s \geq 2$ and that (2.3) holds for all $s' < s$, we get

$$\begin{aligned} \text{sgn}(f \circ h) &= \text{sgn}((f \circ h_1 \circ \dots \circ h_{s-1}) \circ h_s) \\ &= \text{sgn}(f \circ h_1 \circ \dots \circ h_{s-1}) \cdot \text{sgn}(h_s) \\ &= \text{sgn}(f) \cdot \text{sgn}(h_1 \circ \dots \circ h_{s-1}) \cdot \text{sgn}(h_s) \\ &= \text{sgn}(f) \cdot \text{sgn}(h_1 \circ \dots \circ h_{s-1} \circ h_s) \\ &= \text{sgn}(f) \cdot \text{sgn}(h). \end{aligned}$$

Let us note five corollaries of (2.3). First, since $f \circ f^{-1} = \text{id}$, we have $1 = \text{sgn}(\text{id}) = \text{sgn}(f) \cdot \text{sgn}(f^{-1})$ and so $\text{sgn}(f^{-1}) = 1/\text{sgn}(f) = \text{sgn}(f)^{-1}$. Since $1^{-1} = 1$ and $(-1)^{-1} = -1$, we can simplify this formula to read: $\text{sgn}(f^{-1}) = \text{sgn}(f)$ for all $f \in S_n$.

Second, if $h = h_1 \circ \dots \circ h_s$ has been factored as a product of s basic transpositions h_i , then $\text{sgn}(h) = \prod_{i=1}^s \text{sgn}(h_i) = (-1)^s$. For $h = (i, j)$, (2.2) shows that we can write h as the product of an odd number of basic transpositions. Therefore, $\text{sgn}((i, j)) = -1$ for all transpositions (not just basic ones).

Third, if $h = h_1 \circ \dots \circ h_s$ has been factored as a product of s transpositions h_i , then $\text{sgn}(h) = \prod_{i=1}^s \text{sgn}(h_i) = (-1)^s$. In particular, if h can be written in some other way as a product of s' transpositions, then we must have $(-1)^s = (-1)^{s'}$, so that s and s' are both

even or both odd. Thus, while factorizations of a given permutation into transpositions are not unique, the parity of the number of transpositions in such factorizations is unique.

Fourth, if h is a k -cycle, (2.1) shows that h can be written as a product of $k - 1$ transpositions. Therefore, $\text{sgn}(h) = (-1)^{k-1}$. So even-length cycles have negative sign, whereas odd-length cycles have positive sign.

Fifth, suppose $h \in S_n$ is an arbitrary permutation whose disjoint cycle decomposition (including all 1-cycles) has c cycles with respective lengths k_1, \dots, k_c . Then $\text{sgn}(h) = \prod_{i=1}^c (-1)^{k_i-1}$. Since $k_1 + \dots + k_c = n$, we can write $\text{sgn}(h) = (-1)^{n-c}$ where c is the number of cycles in any cycle decomposition of h that includes all 1-cycles.

2.8 Summary

1. *Symmetric groups.* For each positive integer n , let $[n] = \{1, 2, \dots, n\}$. S_n consists of all bijections (permutations) $f : [n] \rightarrow [n]$ with composition of functions as the group operation. S_n is a group of size $n!$, which is non-commutative for $n \geq 3$.
 2. *Directed graphs of functions.* A function $f : [n] \rightarrow [n]$ can be represented by a directed graph with vertex set $[n]$ and a directed edge from i to $f(i)$ for each $i \in [n]$ (which is a loop if $f(i) = i$). For general f , the directed graph of f consists of directed trees feeding into directed cycles. For bijections $f \in S_n$, the directed graph of f consists of one or more directed cycles involving pairwise disjoint sets of vertices. We obtain cycle decompositions of f by traversing all these cycles (with the possible omission of length-1 cycles) and listing the elements encountered.
 3. *Factorizations of permutations.* Every $f \in S_n$ can be factored into a product of disjoint cycles (which commute in S_n); this factorization is unique except for the order of the factors and the possible omission of 1-cycles. Every $f \in S_n$ can be factored as a product of transpositions (i, j) and as a product of basic transpositions $(i, i+1)$. These factorizations are not unique, but the parity (even or odd) of the number of transpositions used is unique.
 4. *Facts about inv.* For any function $f : [n] \rightarrow [m]$, $\text{inv}(f)$ is the number of pairs $i < j$ such that $f(i) > f(j)$. If we sort the list $[f(1), \dots, f(n)]$ by interchanging adjacent elements that are out of order, then we will always make exactly $\text{inv}(f)$ interchanges. Any $f \in S_n$ can be written as the product of $\text{inv}(f)$ basic transpositions.
 5. *Facts about sgn.* For $f \in S_n$, we define $\text{sgn}(f) = (-1)^{\text{inv}(f)}$. For all $f, g \in S_n$ and all $i \neq j$ in $[n]$, $\text{sgn}(f \circ g) = \text{sgn}(f) \cdot \text{sgn}(g)$, $\text{sgn}(\text{id}) = 1$, $\text{sgn}((i, j)) = -1$, $\text{sgn}(g^{-1}) = \text{sgn}(g)$, and $\text{sgn}(f) = (-1)^{n-c}$, where c is the number of cycles in any cycle decomposition of f that includes all 1-cycles. If f can be factored into a product of s transpositions, then $\text{sgn}(f) = (-1)^s$. The sign of any k -cycle is $(-1)^{k-1}$.
-

2.9 Exercises

1. Let $f = [2, 2, 5, 1, 3]$ and $g = [3, 4, 2, 5, 1]$. Find $f \circ g$, $g \circ f$, $f^3 = f \circ f \circ f$, and g^{-1} , giving all answers in one-line form.

2. Let $f = [3, 1, 2, 4]$ and $g = [4, 3, 2, 1]$. Find $f \circ g$, $g \circ f$, $f \circ g \circ f^{-1}$, f^3 , and f^{1001} , giving all answers in one-line form. Explain how you found f^{1001} .
3. Suppose $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ are functions. (a) Prove: if f and g are one-to-one, then $g \circ f$ is one-to-one. (b) Prove: if f and g are onto, then $g \circ f$ is onto. (c) Prove: if f and g are bijections, then $g \circ f$ is a bijection. (d) Prove: if $g \circ f$ is one-to-one, then f is one-to-one. Show by an example that g may not be one-to-one. (e) Prove: if $g \circ f$ is onto, then g is onto. Show by an example that f may not be onto.
4. Prove that function composition is associative; i.e., for all functions $h : W \rightarrow X$, $g : X \rightarrow Y$, and $f : Y \rightarrow Z$, $f \circ (g \circ h) = (f \circ g) \circ h$.
5. Show that the group S_n is not commutative for $n \geq 3$.
6. (a) For $n \geq 1$, explain informally why a function $f : [n] \rightarrow [n]$ is one-to-one iff f is onto. (b) Give an example of a function $f : \mathbb{Z} \rightarrow \mathbb{Z}$ that is one-to-one but not onto. (c) Give an example of a function $g : \mathbb{Z} \rightarrow \mathbb{Z}$ that is onto but not one-to-one.
7. For any set X , let $S(X)$ be the set of all bijections $f : X \rightarrow X$. (a) Prove that $(S(X), \circ)$ is a group. (b) For which sets X is $S(X)$ a commutative group? (c) Suppose $X = \{x_1, \dots, x_n\}$ is an n -element set. Prove that the groups $S(X)$ and S_n are isomorphic.
8. Draw the directed graphs for each of these functions given in one-line form: (a) $[2, 3, 4, 5, 5]$; (b) $[5, 4, 3, 2, 1]$; (c) $[2, 2, 2, 2, 2]$; (d) $[1, 2, 3, 4, 5]$; (e) $[3, 4, 5, 1, 2]$; (f) $[3, 3, 7, 12, 12, 11, 12, 9, 8, 7, 11, 10]$.
9. Let f be the function whose directed graph is drawn in Figure 2.1. Find the one-line forms of $f \circ f$ and $f \circ f \circ f$, and draw the associated directed graphs.
10. Suppose X is an infinite set. Must the directed graph of a bijection $f : X \rightarrow X$ be a disjoint union of directed cycles? If not, can you describe the general structure of this graph?
11. Let $f = [4, 2, 1, 3] \in S_4$. Find **all** possible cycle decompositions of f .
12. Compute each of the following, giving a cycle decomposition for each answer.
 - (a) $(1, 3, 6, 2) \circ (2, 4, 5)(1, 6)$
 - (b) the inverse of $(1, 3, 5, 7, 9)(2)(4, 8, 6)$
 - (c) $[3, 1, 2, 4] \circ (3, 1, 2, 4)$
 - (d) $(3, 1, 2, 4) \circ [3, 1, 2, 4]$
 - (e) $(3, 1, 2, 4) \circ (3, 1, 2, 4)$
 - (f) $(1, 3, 4) \circ (2, 3)(1, 5) \circ (2, 5, 3, 1)^{-1} \circ (2, 3, 4) \circ (1, 2)(3, 4)$
 - (g) $(2, 3) \circ (1, 2) \circ (2, 3) \circ (3, 4) \circ (2, 3) \circ (4, 5) \circ (3, 4) \circ (2, 3) \circ (1, 2)$
13. Let $f = (3, 1, 6)(2, 4)$ and $g = (5, 2, 1, 4, 3)$ in S_6 . (a) Compute $f \circ g$, $g \circ f$, $f \circ g \circ f^{-1}$, and $g \circ f \circ g^{-1}$, giving cycle decompositions for each answer. (b) Compute f^k for $k = 2, 3, 4, 5, 6$. Then find (with explanation) f^{670} .
14. List all $f \in S_4$ with $f \circ f = \text{id}$. Explain why these f 's, and no others, work.
15. (a) Given a cycle decomposition of $f \in S_n$, how can one quickly obtain a cycle decomposition of f^{-1} ? Prove your answer. (b) Given a cycle decomposition of $f \in S_n$ and a large positive integer m , describe how one can find a cycle decomposition of $f^m = f \circ f \circ \dots \circ f$ (m copies of f).
16. (a) Find six permutations $h \in S_8$ such that $h \circ h = (1, 3)(2, 5)(4, 8)(6, 7)$.
 (b) Determine (with proof) how many $h \in S_{4n}$ satisfy

$$h \circ h = (1, 2)(3, 4)(5, 6) \cdots (4n-1, 4n).$$

17. For which $f \in S_n$ does there exist $h \in S_n$ with $h \circ h = f$? Find and prove a necessary and sufficient condition involving the cycle decomposition of f .
18. (a) Give a careful proof of the formula $(b_1, b_2, \dots, b_m) = (b_1, b_2) \circ (b_2, b_3) \circ \dots \circ (b_{m-1}, b_m)$. (b) Give a careful proof of the formula (2.2) in §2.4.
19. *Conjugation rule in S_n .* (a) Suppose $g \in S_n$ and $f = (i_1, i_2, \dots, i_k)$ is a k -cycle in S_n . Prove carefully that $g \circ f \circ g^{-1} = (g(i_1), g(i_2), \dots, g(i_k))$. (b) Prove that the map $C_g : S_n \rightarrow S_n$, given by $C_g(f) = g \circ f \circ g^{-1}$ for $f \in S_n$, is a group isomorphism. (c) Prove that for all $f, g \in S_n$, a cycle decomposition for $g \circ f \circ g^{-1}$ can be obtained from a cycle decomposition for f by applying g to each number in the cycle decomposition, leaving all parentheses unchanged.
20. Find all $f \in S_5$ that commute with $g = (1, 2, 4)(3, 5)$.
21. (a) Suppose the directed graph of $f \in S_n$ consists of k_1 1-cycles, k_2 2-cycles, and so on. Count the number of cycle decompositions of f in which no 1-cycles are omitted and shorter cycles are always listed before longer cycles. (b) Use the conjugation rule (Exercise 19) to prove that the answer to (a) is the number of $g \in S_n$ that commute with f .
22. *Uniqueness of disjoint cycle factorizations.* Let $f \in S_n$ have a cycle decomposition $C_1 C_2 \cdots C_k$ in which all 1-cycles have been omitted. Prove carefully that for any factorization $f = D_1 \circ D_2 \circ \dots \circ D_j$, where the D_i 's are pairwise disjoint non-identity cycles in S_n , we must have $j = k$ and D_1, \dots, D_k is a reordering of C_1, \dots, C_k (viewed as elements of S_n). [Hint: study the cycles of the directed graph of $D_1 \circ \dots \circ D_j$.]
23. Make a table listing all permutations $f \in S_4$. In column 1, write f in one-line form. In column 2, draw the directed graph for f . In column 3, give a cycle decomposition of f . In column 4, compute $\text{inv}(f)$. In column 5, compute $\text{sgn}(f)$.
24. (a) Let $g = [2, 4, 7, 1, 5, 3, 6] \in S_7$. Compute $\text{Inv}(g)$, $\text{inv}(g)$, $\text{sgn}(g)$, $\text{Inv}(g^{-1})$, $\text{inv}(g^{-1})$, and $\text{sgn}(g^{-1})$. (b) Let $h = (2, 4, 7)(1, 5)(3, 6) \in S_7$. Compute $\text{Inv}(h)$, $\text{inv}(h)$, $\text{sgn}(h)$, $\text{Inv}(h^{-1})$, $\text{inv}(h^{-1})$, and $\text{sgn}(h^{-1})$.
25. Suppose $f \in S_{12}$ has a cycle decomposition $f = (a, b, c)(d, e, g, h)(j, k)$. What is $\text{sgn}(f)$?
26. Find and prove a simple relation between $\text{inv}(f)$ and $\text{inv}(f^{-1})$, for $f \in S_n$. Deduce that $\text{sgn}(f) = \text{sgn}(f^{-1})$.
27. Write each permutation as a product of basic transpositions. (a) $f = [4, 1, 3, 5, 2]$; (b) $g = (3, 6)$; (c) $h = (1, 2, 4)(3, 5)$.
28. Find, with explanation: (a) the number of basic transpositions in S_n ; (b) the number of transpositions in S_n ; (c) the number of n -cycles in S_n .
29. For even $n \geq 2$, how many $f \in S_n$ satisfy $f(1) < f(3) < f(5) < \dots < f(n-1)$ and $f(i) < f(i+1)$ for all odd $i \in [n]$?
30. (a) Write $f = [3, 7, 4, 1, 6, 8, 2, 5]$ as a product of transpositions in three different ways. (b) What is the least number of transpositions appearing in any such factorization for f ?
31. Find all $f \in S_5$ with $\text{inv}(f) = 8$.
32. Verify assertions (1) and (2) in §2.6.
33. Recall that we defined $\text{sgn}(f) = 0$ when f is not a bijection. With this convention, prove that for all functions $f, g : [n] \rightarrow [n]$, $\text{sgn}(f \circ g) = \text{sgn}(f) \cdot \text{sgn}(g)$.

34. Fix $1 \leq i < j \leq n$. Use the definitions to compute $\text{Inv}((i, j))$, $\text{inv}((i, j))$, and $\text{sgn}((i, j))$.
35. Which permutation in S_n has the most inversions? Give the one-line form and a cycle decomposition for this permutation.
36. (a) Suppose $f : [n] \rightarrow [n]$ has one-line form $[f(1), \dots, f(n)]$, and g is obtained from this one-line form by switching $f(i)$ and $f(j)$ (the entries in *positions* i and j). Prove carefully that $g = f \circ (i, j)$. (b) Suppose $h = (i, j) \circ f$. Describe (with proof) how the one-line form for h can be obtained from the one-line form for f . (c) Use (a) or (b) and an argument involving sorting the one-line form for f to prove that any $f \in S_n$ can be written as a product of transpositions.
37. (a) For $n \geq 2$, show that $A_n = \{f \in S_n : \text{sgn}(f) = +1\}$ is a normal subgroup of S_n . (b) Show that $V = \{(1), (1, 2)(3, 4), (1, 3)(2, 4), (1, 4)(2, 3)\}$ is a normal subgroup of S_4 . (c) Show that $H = \{(1), (1, 2, 3), (1, 3, 2)\}$ is a subgroup of S_4 that is not normal in S_4 .
38. Prove: for all $n \geq 3$ and all $f \in S_n$ with $\text{sgn}(f) = +1$, f can be factored into a product of (not necessarily disjoint) 3-cycles.
39. Prove: for $n \geq 2$, exactly half of the permutations in S_n have sign -1 . [Hint: consider the map sending each $f \in S_n$ to $f \circ (1, 2)$.]
40. (a) Prove that any $f \in S_n$ can be written as a product of factors, each of which is either $(1, 2)$ or $(1, 2, \dots, n)$. Illustrate your proof with $f = [3, 5, 1, 4, 2]$. (b) What is the relation between the one-line forms of f and $f \circ (1, 2)$? What is the relation between the one-line forms of f and $f \circ (1, 2, \dots, n)$? (c) Interpret parts (a) and (b) as a result about the ability to sort a list into increasing order using just two particular sorting moves. (d) What can you say about the minimum number of steps needed to sort a list using the two moves in part (c)? What if right-multiplication by $(n, \dots, 2, 1)$ is also allowed as a move?
41. True or false? Explain each answer. (a) $(3, 1, 2, 4)$ is a 3-cycle. (b) $[3, 1, 2, 4]$ is a 3-cycle. (c) There exists $f \in S_4$ with $[f(1), f(2), f(3), f(4)] = (f(1), f(2), f(3), f(4))$. (d) For all functions $f : [n] \rightarrow [n]$, if $f \circ f \in S_n$ then $f \in S_n$. (e) For all $f \in S_n$, if $f^{-1} = f$ then $f = \text{id}$ or f is a 2-cycle. (f) For all distinct cycles $f, g \in S_n$, $f \circ g = g \circ f$ iff f and g are disjoint cycles. (g) Any product of c disjoint cycles in S_n has sign $(-1)^{n-c}$. (h) For all $1 < k \leq n$, every k -cycle in S_n can be written as a product of $k - 1$ basic transpositions. (i) For all $n \geq 2$ and all $f \in S_n$ and all positive integers k , there is a way to write f as a product of $\text{inv}(f) + 2k$ transpositions. (j) For all $f, g \in S_n$ and $1 \leq i < n$, if $g = (i, i + 1) \circ f$, then $\text{inv}(g) = \text{inv}(f) + 1$ or $\text{inv}(g) = \text{inv}(f) - 1$. (k) Every $g \in S_n$ can be written as a product of disjoint transpositions. (l) For all $f \in S_n$, if $f = h \circ h$ for some $h \in S_n$, then $\text{sgn}(f) = +1$. (m) For all $f \in S_n$, if $\text{sgn}(f) = +1$ then $f = h \circ h$ for some $h \in S_n$.

This page intentionally left blank

3

Polynomials

Polynomials appear ubiquitously in linear algebra and throughout mathematics. Most of us encounter the idea of a polynomial *function* in calculus. In algebraic settings, one needs a more formal concept of polynomial as a symbolic expression of the form $\sum_{i=0}^n a_i x^i$ that is not necessarily a “function of x .” We begin this chapter with an intuitive description of these formal polynomials, which form a ring under the familiar algebraic rules for adding and multiplying polynomials. Then we make this intuitive discussion more precise by giving a rigorous definition of polynomials in terms of formal power series and connecting this definition with the idea of a polynomial function.

The next part of the chapter explores the divisibility structure of one-variable polynomials with coefficients in a field. We will see a close analogy between factorization of such polynomials and more elementary results on factorization of integers. Recall that for any integers a, b with $b \neq 0$, there is a *division algorithm* that produces a unique quotient $q \in \mathbb{Z}$ and remainder $r \in \mathbb{Z}$ satisfying $a = bq + r$ and $0 \leq r < |b|$. We will develop a similar division algorithm for polynomials, where the remainder must now have smaller degree than b . This algorithm will elucidate properties of divisors, greatest common divisors, and least common multiples of polynomials. To understand how polynomials factor, we introduce *irreducible* polynomials, which are non-constant polynomials that cannot be factored into products of two other polynomials of smaller degree. Irreducible polynomials are the analogues of prime integers. We will see that every nonzero polynomial has a unique factorization into a product of irreducible polynomials, just as every nonzero integer can be uniquely written as a product of primes.

The chapter concludes with a deeper study of irreducible polynomials. We present some theorems and algorithms for testing whether polynomials with coefficients in certain fields are irreducible. In particular, Kronecker’s factoring algorithm gives a laborious but definitive method for testing a polynomial with rational coefficients for irreducibility, or finding the irreducible factorization of such a polynomial. We also give a brief discussion of minimal polynomials, which are polynomials that help us understand the structure of finite-dimensional F -algebras.

3.1 Intuitive Definition of Polynomials

We begin by presenting an informal definition of polynomials, which will be made more precise shortly. Let R be any ring. Intuitively, a *one-variable polynomial with coefficients in R* is some expression of the form

$$a_0 + a_1 x + a_2 x^2 + \cdots + a_n x^n$$

where $n \geq 0$ is an integer, each a_i is an element of R , and the letter x is an “indeterminate” or “formal variable.” Often we denote such a polynomial using the summation notation $\sum_{i=0}^n a_i x^i$. It is convenient to define a_i for all integers $i \geq 0$ by setting $a_i = 0_R$ for $i > n$.

Then, as a matter of definition, two polynomials $\sum_{i \geq 0} a_i x^i$ and $\sum_{i \geq 0} b_i x^i$ are *equal* iff $a_i = b_i$ for all $i \geq 0$. So, two polynomials are the same iff all their coefficients agree. A polynomial is only permitted to have finitely many nonzero coefficients. The symbol $R[x]$ denotes the set of all polynomials with coefficients in R . For example, $1 + 3x - \pi x^2 + \sqrt{7}x^4 \in \mathbb{R}[x]$, but $\sum_{n=0}^{\infty} (1/n!)x^n$ is not a polynomial.

The unique polynomial whose coefficients are all zero is called the *zero polynomial*. The *degree* of a nonzero polynomial $p = \sum_{i \geq 0} a_i x^i$ is the largest integer n such that $a_n \neq 0_R$. The degree of the zero polynomial is undefined. If p is a polynomial of degree n , we write $\deg(p) = n$ and introduce the following terminology: $a_n x^n$ is called the *leading term* of p ; x^n is called the *leading monomial* of p ; and a_n is called the *leading coefficient* of p . We call p a *monic* polynomial iff $a_n = 1$. Polynomials of degree zero, one, two, three, four, and five are respectively called *constant polynomials*, *linear polynomials*, *quadratic polynomials*, *cubic polynomials*, *quartic polynomials*, and *quintic polynomials*. The zero polynomial is also considered to be a constant polynomial.

3.2 Algebraic Operations on Polynomials

Given any ring R , we can turn $R[x]$ into a ring by introducing the following addition and multiplication operations on polynomials. Suppose $p = \sum_{i \geq 0} a_i x^i$ and $q = \sum_{j \geq 0} b_j x^j$ are two elements of $R[x]$, where $a_i, b_j \in R$. We define the *sum* of p and q by

$$p + q = \sum_{k \geq 0} (a_k + b_k)x^k.$$

The distributive law suggests that the product pq should be given by

$$pq = \sum_{i \geq 0} \sum_{j \geq 0} a_i b_j x^{i+j}.$$

To present this in standard form, we need to collect like powers of x . In the preceding expression, the coefficient of x^k is the sum of all products $a_i b_j$ such that $i + j = k$. So we define the *product* of p and q by

$$pq = \sum_{k \geq 0} \left(\sum_{i+j=k} a_i b_j \right) x^k = \sum_{k \geq 0} \left(\sum_{i=0}^k a_i b_{k-i} \right) x^k.$$

(In summations like this, it is always understood that the summation variables only take nonnegative integer values.) With these definitions, it is now tedious but straightforward to check that the ring axioms do hold. One must confirm that addition of polynomials is closed, associative, commutative, has an identity element (the zero polynomial), and has additive inverses. Similarly, multiplication of polynomials is closed, associative, has an identity element (the constant polynomial $1_{R[x]} = 1_R + 0_R x + 0_R x^2 + \dots$), and distributes over addition. When the coefficient ring R is *commutative*, $R[x]$ is also commutative. We verify a few of the ring axioms now, leaving the others as exercises.

First, consider the additive inverse axiom. Assume we already have checked that the zero polynomial is the additive identity of $R[x]$. Given $p = \sum_{i=0}^n a_i x^i \in R[x]$, does p have an additive inverse in $R[x]$? Since R is a ring, each coefficient $a_i \in R$ has an additive inverse $-a_i \in R$, so that $q = \sum_{i=0}^n (-a_i) x^i$ is an element of $R[x]$. By definition of addition of

polynomials, we compute

$$p + q = \sum_{i=0}^n (a_i + (-a_i))x^i = \sum_{i=0}^n 0_R x^i = 0_{R[x]} = \sum_{i=0}^n ((-a_i) + a_i)x^i = q + p.$$

Therefore q is an additive inverse for p in $R[x]$.

Second, consider associativity of multiplication. Given $p, q, r \in R[x]$, write $p = \sum_{i \geq 0} a_i x^i$, $q = \sum_{j \geq 0} b_j x^j$, and $r = \sum_{k \geq 0} c_k x^k$, where all $a_i, b_j, c_k \in R$. We have $pq = \sum_{t \geq 0} d_t x^t$ where $d_t = \sum_{i+j=t} a_i b_j$, and $qr = \sum_{u \geq 0} e_u x^u$ where $e_u = \sum_{j+k=u} b_j c_k$. Applying the definition of product again, we see that $(pq)r = \sum_{v \geq 0} f_v x^v$ where

$$f_v = \sum_{t+k=v} d_t c_k = \sum_{(t,k): t+k=v} \left(\sum_{(i,j): i+j=t} a_i b_j \right) c_k = \sum_{(i,j,k): (i+j)+k=v} (a_i b_j) c_k.$$

The last step used the distributive law in the ring R , as well as commutativity and associativity of addition in R . Likewise, $p(qr) = \sum_{v \geq 0} g_v x^v$ where

$$g_v = \sum_{i+u=v} a_i e_u = \sum_{(i,u): i+u=v} a_i \left(\sum_{(j,k): j+k=u} b_j c_k \right) = \sum_{(i,j,k): i+(j+k)=v} a_i (b_j c_k).$$

Since $(a_i b_j) c_k = a_i (b_j c_k)$ holds for all i, j, k by associativity in R , we see that $f_v = g_v$ for all $v \geq 0$. Therefore, by definition of equality of polynomials, $(pq)r = p(qr)$.

Now assume that F is a field. We can define a scalar multiplication on polynomials in $F[x]$ by setting $c \cdot (\sum_{i=0}^n a_i x^i) = \sum_{i=0}^n (ca_i)x^i$ for $c, a_i \in F$. Routine verifications show that $F[x]$, under polynomial addition and this scalar multiplication, is a vector space over F . Moreover, $B = \{1, x, x^2, x^3, \dots, x^n, \dots\}$ is a basis for $F[x]$, which is therefore an infinite-dimensional F -vector space. For $c \in F$ and $p, q \in F[x]$, one sees that $c \cdot (pq) = (c \cdot p)q = p(c \cdot q)$, so that $F[x]$ is in fact an F -algebra (see §1.3). In the case of an arbitrary ring R , we see similarly that $R[x]$ is an R -module.

3.3 Formal Power Series and Polynomials

Initially, we defined a polynomial $p \in R[x]$ to be an “expression” of the form $a_0 + a_1 x + \dots + a_n x^n$, with each $a_i \in R$. However, we were rather vague about the meaning of the letter x in this expression. In order to build a sound theory, we need to have a more rigorous definition of what a polynomial is. The key to the formal definition of polynomials is our earlier observation that two polynomials should be equal iff they have the same sequence of coefficients. This suggests *identifying* a polynomial $\sum_{i=0}^n a_i x^i$ with its sequence of coefficients (a_0, a_1, \dots, a_n) . For many purposes (e.g., adding two polynomials of different degrees), it is technically more convenient to use an infinite coefficient sequence $(a_0, a_1, \dots, a_n, 0_R, 0_R, \dots)$ that ends with an infinite string of zero coefficients.

Based on this idea of a coefficient sequence, we define a *formal power series with coefficients in the ring R* to be an infinite sequence $f = (f_0, f_1, f_2, \dots, f_i, \dots) = (f_i : i \geq 0)$ with every $f_i \in R$. Two formal power series $f = (f_i : i \geq 0)$ and $g = (g_i : i \geq 0)$ are *equal* iff $f_i = g_i$ for all $i \geq 0$; this is the ordinary meaning of equality of two sequences. A formal power series f is called a *formal polynomial* iff there exists $N \in \mathbb{N}$ with $f_n = 0_R$ for all $n > N$ (i.e., f ends in an infinite string of zeroes). We write $R[[x]]$ to denote the set of all

formal power series with coefficients in R , and (as above) we write $R[x]$ to denote the subset of $R[[x]]$ consisting of the formal polynomials. By analogy with our intuitive notation for polynomials, we often write $f = \sum_{i \geq 0} f_i x^i = \sum_{i=0}^{\infty} f_i x^i$ to denote a formal power series, and we call f_i “the coefficient of x^i in f . ” However, it should be understood that (for now) the summation symbol here is only a notational device — no summation (finite or infinite) is being performed! Similarly, the powers x^i appearing here are (for now) only symbols that form part of the notation for a formal power series or a formal polynomial.

Extending the earlier formulas for polynomial addition and multiplication to formal power series, the following definitions suggest themselves. Suppose R is a ring and $f = (f_i : i \geq 0)$, $g = (g_i : i \geq 0)$ are two elements of $R[[x]]$. Define $f + g = (f_i + g_i : i \geq 0)$. Define $fg = (h_k : k \geq 0)$, where $h_k = \sum_{i+j=k} f_i g_j = \sum_{i=0}^k f_i g_{k-i}$ for each $k \geq 0$. Note that any particular coefficient in $f + g$ or fg can be found by computing *finitely* many sums and products in the underlying ring R ; so our definitions of $f + g$ and fg do not rely on limits or any other infinite process. As in the case of polynomials, we can verify that $R[[x]]$ with these operations is a ring, which is commutative if R is commutative. For example, the proofs of the additive inverse axiom and associativity of multiplication, which we gave earlier for polynomials, carry over almost verbatim to the present situation. As another example, let us see that $R[[x]]$ is closed under multiplication. Suppose $f, g \in R[[x]]$, and write $f = (f_i : i \geq 0)$, $g = (g_i : i \geq 0)$ for some $f_i, g_i \in R$. By definition, $fg = (h_k : k \geq 0)$ where $h_k = \sum_{i=0}^k f_i g_{k-i}$. To see that $fg \in R[[x]]$, we must prove that $h_k \in R$ for all $k \geq 0$. For fixed $k \geq 0$ and $0 \leq i \leq k$, $f_i g_{k-i}$ is in R because $f_i \in R$, $g_{k-i} \in R$, and R is closed under multiplication. Since h_k is a sum of finitely many elements of this form, h_k is also in R because R is closed under addition. So $fg \in R[[x]]$.

Once we have proved all the ring axioms for $R[[x]]$, we can reprove that $R[x]$ is a ring by verifying that $R[x]$ is a *subring* of $R[[x]]$. Since the formal series $0 = (0, 0, 0, \dots)$ and $1 = (1, 0, 0, \dots)$ are polynomials, it suffices to check that for all $f, g \in R[x]$, $f + g \in R[x]$ and $-f \in R[x]$ and $fg \in R[x]$. In other words, if the sequences f and g each end in infinitely many zeroes, the same holds for $f + g$ and $-f$ and fg . We let the reader check this assertion in the case of $f + g$ and $-f$. To treat the case of fg , fix $N, M \in \mathbb{N}$ such that $f_n = 0$ for all $n > N$ and $g_m = 0$ for all $m > M$. For any fixed $k > N + M$, the coefficient of x^k in fg is $\sum_{n+m=k} f_n g_m$. Since $k > N + M$, the condition $n + m = k$ forces $n > N$ or $m > M$, so $f_n = 0$ or $g_m = 0$, and $f_n g_m = 0$. Thus, the coefficient in question is a sum of zeroes, so is 0.

We can now give a rigorous meaning to our original notation for polynomials as “expressions” involving x . Define the letter x to be the particular sequence $(0, 1, 0, 0, \dots)$, which is both a polynomial in $R[x]$ and a formal series in $R[[x]]$. This definition is intuitively reasonable upon recalling that this coefficient sequence was designed as the formal model for the informal polynomial expression $0 + 1x^1 + 0x^2 + 0x^3 + \dots$. Iterating the definition of multiplication, one checks by induction that for all $i \geq 0$, x^i (i.e., the product of i copies of the sequence x in $R[[x]]$) is the sequence $(e_j : j \geq 0)$, where $e_i = 1$ and $e_j = 0$ for all $j \neq i$. Next, we can identify any ring element $a \in R$ with the constant polynomial $(a, 0, 0, \dots) \in R[[x]]$. More precisely, the map $a \mapsto (a, 0, 0, \dots)$ is an injective ring homomorphism of R into $R[[x]]$. We use this injection to regard R as a subset of $R[[x]]$ (or $R[x]$), writing $a = (a, 0, 0, \dots)$. Finally, one can check that

$$a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n = (a_0, a_1, a_2, \dots, a_n, 0, 0, \dots),$$

where the left side is built up by performing sum and product operations in the *ring* $R[x]$. (Similarly, one can justify the notation $\sum_{i=0}^{\infty} f_i x^i$ for a formal power series as actually encoding a sum of terms, each of which is a product of f_i and i copies of x , but here one needs to introduce a limiting device of some kind to define what the infinite sum means.)

For a field F , $F[[x]]$ is an F -vector space and F -algebra with scalar multiplication defined by $c \cdot (f_i : i \geq 0) = (cf_i : i \geq 0)$ for $c, f_i \in F$. $F[x]$ is a subspace of $F[[x]]$. (Likewise, for a ring R , $R[[x]]$ is an R -module with submodule $R[x]$.) Using the explicit description of the sequences x^i given in the previous paragraph, one readily checks that the set $B = \{1, x, x^2, \dots, x^n, \dots\}$ is an F -linearly independent subset of $F[[x]]$. However, B is *not* a basis of $F[[x]]$. The reason is that the span of B consists of all *finite* F -linear combinations of the powers of x , which are precisely the elements of $F[x]$. The formal series $(1_F, 1_F, 1_F, \dots) = \sum_{n=0}^{\infty} x^n$ is an example of an element of $F[[x]]$ outside the span of B . It is true that B can be extended to a basis of $F[[x]]$, but it does not seem possible to give an explicit description of any specific basis of $F[[x]]$.

3.4 Properties of Degree

In this section, we study how the algebraic operations in $R[x]$ affect degree. Let R be any ring, and suppose p, q are nonzero polynomials in $R[x]$. Write $p = \sum_{i=0}^N a_i x^i$ and $q = \sum_{j=0}^M b_j x^j$ where $N = \deg(p)$, $M = \deg(q)$, and $a_i, b_j \in R$. By definition of addition, we see that $\deg(p + q) = N$ if $N > M$, whereas $\deg(p + q) = M$ if $M > N$. On the other hand, when $N = M$, it is possible for the leading terms to cancel. More precisely, when $N = M$, we have: $\deg(p + q) = N$ if $b_N \neq -a_N$; $\deg(p + q) < N$ if $b_N = -a_N$ and $p + q \neq 0$; and $\deg(p + q)$ is undefined if $p + q = 0$. To summarize, we always have

$$\deg(p + q) \leq \max(\deg(p), \deg(q))$$

provided that p, q , and $p + q$ are nonzero, and strict inequality only occurs when $\deg(p) = \deg(q)$ and the leading terms of p and q cancel.

Next consider the degree of the product pq , where $\deg(p) = N$ and $\deg(q) = M$ as above. Earlier, as part of the proof that $R[x]$ is a subring of $R[[x]]$, we showed that the coefficient of x^k in pq is zero for all $k > N + M$. We use a similar argument to find the coefficient of x^{N+M} in pq . By definition, this coefficient is $\sum_{i+j=N+M} a_i b_j$. One term in this sum is $a_N b_M$. For all the other terms, either $i > N$ and $a_i = 0$, or $j > M$ and $b_j = 0$. It follows that

$$pq = a_N b_M x^{N+M} + \text{terms involving lower powers of } x.$$

We can deduce several consequences from this computation. First, *for all nonzero $p, q \in R[x]$, if the leading terms of p and q do not multiply to zero, then*

$$pq \neq 0 \text{ and } \deg(pq) = \deg(p) + \deg(q).$$

Let us call the displayed conclusion the *degree addition formula*. Note that the degree addition formula always holds if R is an integral domain (which has no zero divisors) or a field (which is a special kind of integral domain). In fact, R is an integral domain iff $R[x]$ is an integral domain. The forward implication is part of the degree addition formula ($p, q \neq 0$ implies $pq \neq 0$), and the converse follows since R is isomorphic to a subring of $R[x]$, namely the subring of constant polynomials.

On the other hand, one can check that for any ring R , $R[x]$ is never a field. In the particular case of an integral domain R , we can show that $p \in R[x]$ has a multiplicative inverse in $R[x]$ iff $\deg(p) = 0$ (i.e., p is a nonzero constant polynomial) and p is invertible viewed as an element of R . For, assuming $p \in R[x]$ has inverse $q \in R[x]$, applying the degree addition formula to $pq = 1$ shows $\deg(p) + \deg(q) = \deg(1) = 0$. Since degrees are

nonnegative integers, this forces $\deg(p) = \deg(q) = 0$, and from this we also see that q (being a constant polynomial) can be regarded as an inverse of p in the ring R . The converse is immediate, since if $pq = 1$ holds for given $p, q \in R$, then the same relation holds between p, q regarded as constant polynomials in $R[x]$. In the case of a field F , the result just proved is that the invertible elements of $F[x]$ are precisely the nonzero constant polynomials.

Given a field F , we can use degree to define some finite-dimensional subspaces of the F -vector space $F[x]$. Specifically, for each $n \geq 0$, let V_n be the set of all $p \in F[x]$ of degree at most n , together with the zero polynomial. One readily checks that each V_n is an $(n + 1)$ -dimensional subspace of $F[x]$ with ordered basis $(1, x, x^2, \dots, x^n)$.

3.5 Evaluating Polynomials

From the outset, we have stressed that a formal polynomial is a sequence of coefficients, not a function of the “variable” x . Nevertheless, we can use a formal polynomial to create a polynomial function, as follows. Given a ring R and a polynomial $p = \sum_{i=0}^n a_i x^i \in R[x]$, define the *polynomial function* $f_p : R \rightarrow R$ associated to p by setting $f_p(c) = \sum_{i=0}^n a_i c^i$ for each $c \in R$. Note that the expression defining $f_p(c)$ is a formula involving the addition and multiplication operations in the given ring R . We call $f_p(c)$ the *value of the polynomial p at the point $x = c$* . One often writes $p(c)$ instead of $f_p(c)$, although the latter notation is more precise.

One must take care to distinguish the (formal) polynomial p from its associated polynomial function f_p . We reiterate that p is the sequence $(a_0, a_1, \dots, a_n, 0, 0, \dots)$, while f_p is a certain function from R to R . This distinction may at first seem pedantic or trivial, but it is essential. Indeed, it is quite possible for two different polynomials to induce the *same* function from R to R . For example, let $R = \mathbb{Z}_2 = \{0, 1\}$ be the ring of integers modulo 2. There are only four distinct functions $g : \mathbb{Z}_2 \rightarrow \mathbb{Z}_2$, but there are infinitely many polynomials $p \in \mathbb{Z}_2[x]$ (which are sequences of zeroes and ones terminating in an infinite string of zeroes). Thus, one of these four functions must arise from infinitely many distinct polynomials. For an explicit example, consider the function $g : \mathbb{Z}_2 \rightarrow \mathbb{Z}_2$ such that $g(0) = 0$ and $g(1) = 1$. For any $i \geq 1$, the polynomial $x^i \in \mathbb{Z}_2[x]$ has g as its associated polynomial function, since $f_{x^i}(0) = 0^i = 0 = g(0)$ and $f_{x^i}(1) = 1^i = 1 = g(1)$. However, the polynomials $x^1 = (0, 1, 0, 0, \dots)$, $x^2 = (0, 0, 1, 0, 0, \dots)$, etc., are all *distinct* elements of $\mathbb{Z}_2[x]$.

We obtained the function f_p from a given polynomial $p = \sum_{i \geq 0} a_i x^i$ by regarding x as a “variable” that ranges over the elements of the coefficient ring R . In the further study of polynomial rings, it is often helpful to reverse this viewpoint, holding the value of x fixed and letting p range over all polynomials in $R[x]$. More precisely, suppose we are given a fixed ring element $c \in R$. Then we obtain a map $E_c : R[x] \rightarrow R$ defined by $E_c(p) = f_p(c) = p(c)$ for $p \in R[x]$. In words, the map E_c sends each polynomial in $R[x]$ to the value of that polynomial at $x = c$. We call E_c *evaluation at $x = c$* . Assuming R is commutative, one may check that every map E_c is a *ring homomorphism*; this means that $(p + q)(c) = p(c) + q(c)$ and $(pq)(c) = p(c)q(c)$ for all $p, q \in R[x]$ and all $c \in R$, and $E_c(1_{R[x]}) = 1_R$. (The verification that E_c preserves multiplication requires commutativity of R .)

The fact that E_c is a ring homomorphism is a special case of the following more general construction, which is called the *universal mapping property (UMP)¹ of polynomial rings*.

¹Universal mapping properties appear throughout algebra and provide a unifying framework for understanding various algebraic constructions. Chapter 19 gives a systematic introduction to universal mapping properties.

Suppose R and S are commutative rings, and $h : R \rightarrow S$ is a ring homomorphism. For every $c \in S$, there exists a unique ring homomorphism $H : R[x] \rightarrow S$ such that $H(x) = c$ and $H(r) = h(r)$ for all $r \in R$ (i.e., H extends h). We first prove the uniqueness of H . Suppose $H : R[x] \rightarrow S$ is a ring homomorphism extending h such that $H(x) = c$. Let $p = \sum_{i=0}^n a_i x^i$ be an arbitrary polynomial in $R[x]$. As explained near the end of §3.3, the displayed expression for p can be regarded as a sum in $R[x]$ of certain products of polynomials in $R[x]$ (namely a_i times several copies of x). Since H is a ring homomorphism, we must have

$$H(p) = \sum_{i=0}^n H(a_i x^i) = \sum_{i=0}^n H(a_i) H(x)^i = \sum_{i=0}^n h(a_i) c^i.$$

Thus the value of $H(p)$ is completely determined by h and c , proving uniqueness.

To prove existence, use the formula just displayed to define a function $H : R[x] \rightarrow S$. Identifying each $a \in R$ with the constant polynomial $(a, 0, 0, \dots)$, we see that $H(a) = h(a)c^0 = h(a)$ for all $a \in R$, so that H does extend h . In particular, $H(1_{R[x]}) = h(1_R) = 1_S$. Recalling that $x = (0, 1, 0, 0, \dots)$, we see similarly that $H(x) = h(1_R)c^1 = c$. All that remains is to check that H preserves sums and products. Let $p = \sum_{i \geq 0} a_i x^i$ and $q = \sum_{i \geq 0} b_i x^i$ be two polynomials in $R[x]$, where $a_i, b_i \in R$. By definition,

$$H(p+q) = \sum_{i \geq 0} h(a_i + b_i) c^i = \sum_{i \geq 0} (h(a_i) + h(b_i)) c^i = \sum_{i \geq 0} h(a_i) c^i + \sum_{i \geq 0} h(b_i) c^i = H(p) + H(q).$$

Next, since $pq = \sum_{i \geq 0} (\sum_{u+v=i} a_u b_v) x^i$, we find that

$$\begin{aligned} H(pq) &= \sum_{i \geq 0} h\left(\sum_{u+v=i} a_u b_v\right) c^i = \sum_{i \geq 0} \left(\sum_{u+v=i} h(a_u) h(b_v)\right) c^i \\ &= \sum_{i \geq 0} \sum_{u+v=i} h(a_u) c^u h(b_v) c^v \quad (\text{since } S \text{ is commutative and } c^u c^v = c^i) \\ &= \left(\sum_{u \geq 0} h(a_u) c^u\right) \left(\sum_{v \geq 0} h(b_v) c^v\right) \quad (\text{using the distributive law in } S) \\ &= H(p)H(q). \end{aligned}$$

This completes the verification of the UMP. To obtain the evaluation maps E_c mentioned earlier, take $S = R$ and $h = \text{id}_R$.

The UMP is frequently applied to the situation where S is a ring containing R as a subring, and $h : R \rightarrow S$ is the inclusion map. In this case, for every $c \in S$, we obtain a ring homomorphism $H_c : R[x] \rightarrow S$ such that, for each $p = \sum_{i \geq 0} a_i x^i$ in $R[x]$, $H_c(p) = p(c) = \sum_{i \geq 0} a_i c^i$. The image of this ring homomorphism is denoted $R[\{c\}]$ or $R[c]$. This image is a subring of S , and it is the smallest subring of S containing both R and c .

3.6 Polynomial Division with Remainder

Our next topic is the theory of divisibility in polynomial rings. The cornerstone of this theory is the following *division theorem*, which precisely describes the process of polynomial long division with remainder. *Given any $f, g \in F[x]$, where F is a field and g is nonzero, there exist unique polynomials $q, r \in F[x]$ such that $f = gq+r$ and either $r = 0$ or $\deg(r) < \deg(g)$.* We call q the *quotient* and r the *remainder* when f is divided by g .

First we prove the existence of q and r . If $f = 0$, we may take $q = r = 0$. If f is nonzero, we may proceed by strong induction on $n = \deg(f)$. The base case occurs when $n < \deg(g)$; here we may take $q = 0$ and $r = f$. For the induction step, fix a polynomial f with $\deg(f) = n \geq \deg(g)$, and assume that the existence assertion is already known for all polynomials of smaller degree than f . Write $f = ax^n +$ (lower terms) and $g = bx^m +$ (lower terms), where a and b are the leading coefficients of f and g , and $n \geq m$. Since b is a nonzero element of the field F , its inverse b^{-1} exists in F . Consider now the polynomial $h = f - (ab^{-1}x^{n-m})g \in F[x]$. The polynomial h is the difference of two polynomials, both of whose leading terms equal ax^n . It follows that $h = 0$ or $\deg(h) < n$. If $h = 0$, we may take $q = ab^{-1}x^{n-m}$ and $r = 0$. Otherwise, the induction hypothesis is applicable to h and yields polynomials $q_0, r_0 \in F[x]$ such that $h = q_0g + r_0$ and either $r_0 = 0$ or $\deg(r_0) < \deg(g)$. Now, $f = qg + r$ holds if we take $q = ab^{-1}x^{n-m} + q_0$ and $r = r_0$. Since we chose $r = r_0$, we have $r = 0$ or $\deg(r) < \deg(g)$.

Let us now prove that the quotient and remainder associated with a given f and g are unique. Suppose we had $f = q_1g + r_1$ and $f = q_2g + r_2$ where $q_1, q_2, r_1, r_2 \in F[x]$, $r_1 = 0$ or $\deg(r_1) < \deg(g)$, and $r_2 = 0$ or $\deg(r_2) < \deg(g)$. We need to prove that $q_1 = q_2$ and $r_1 = r_2$. Equating the given expressions for f and rearranging, we find that $(q_1 - q_2)g = r_2 - r_1$. The polynomial $r_2 - r_1$ is either zero or nonzero. In the case where $r_2 - r_1 = 0$, we deduce $r_1 = r_2$ and $(q_1 - q_2)g = 0$. Since F is a field, $F[x]$ is an integral domain. Then, since $g \neq 0$, it follows that $q_1 - q_2 = 0$ and $q_1 = q_2$. Thus the uniqueness result holds. In the other case, where $r_2 - r_1 \neq 0$, it follows that $q_1 - q_2 \neq 0$. So we may compute degrees, obtaining

$$\deg(r_2 - r_1) = \deg((q_1 - q_2)g) = \deg(q_1 - q_2) + \deg(g) \geq \deg(g).$$

(This calculation uses the degree addition formula in $F[x]$, which is valid because the field F is an integral domain.) But the assumptions on r_1 and r_2 ensure that $r_2 - r_1$ has degree strictly less than $\deg(g)$. This contradiction shows that the second case never occurs.

We close this section with two further remarks, one computational and one theoretical. First, observe that the induction step in the existence proof can be unravelled into the familiar iterative algorithm for long division of polynomials learned in elementary algebra. For example, suppose we are dividing $f = 3x^3 - 7x^2 + 4$ by $g = 6x + 1$ in the ring $\mathbb{Q}[x]$. The first step in the long division is to divide the leading terms of f and g , obtaining $3x^3/6x = (1/2)x^2$. This monomial corresponds to the expression $ab^{-1}x^{n-m}$ in the proof above. As prescribed in the proof, we now multiply all of g by this monomial, subtract the result from f , and continue the division process. Usually the computations are presented in a tabular form as shown in Figure 3.1. The process ends when the “leftover part” is either zero or has degree less than the degree of g . In this example, once the division is complete, we read off the quotient $q = (1/2)x^2 - (5/4)x + 5/24$ and the remainder $r = 91/24$.

The second remark concerns the possibility of extending the division theorem to polynomial rings $R[x]$ where R may not be a field. We restrict attention to commutative rings R . For convenience, we have italicized the steps in the proof above that used the assumption that F was a field. We see that the uniqueness proof will go through provided that $R[x]$ is an integral domain, which (as we have seen) holds iff R itself is an integral domain. On the other hand, to make the existence proof work, we needed the leading coefficient of g (called b in the proof) to be an invertible element of the coefficient ring. We therefore obtain the following generalization of the division theorem: *Given a commutative ring R and polynomials $f, g \in R[x]$ such that $g \neq 0$ and the leading coefficient of g is invertible in R , there exist $q, r \in R[x]$ with $f = qg + r$ and $r = 0$ or $\deg(r) < \deg(g)$. If R is an integral domain, then q and r are unique.* The new theorem applies, in particular, if we wish to divide an arbitrary polynomial $f \in \mathbb{Z}[x]$ by a monic polynomial $g \in \mathbb{Z}[x]$. The resulting quotient and remainder will lie in $\mathbb{Z}[x]$ and be unique.

$$\begin{array}{r}
 \frac{\frac{1}{2}x^2 - \frac{5}{4}x + \frac{5}{24}}{6x + 1} \\
 \hline
 3x^3 - 7x^2 + 4 \\
 - 3x^3 - \frac{1}{2}x^2 \\
 \hline
 - \frac{15}{2}x^2 \\
 \frac{15}{2}x^2 + \frac{5}{4}x \\
 \hline
 \frac{5}{4}x + 4 \\
 - \frac{5}{4}x - \frac{5}{24} \\
 \hline
 \frac{91}{24}
 \end{array}$$

FIGURE 3.1

Example of Polynomial Division with Remainder.

3.7 Divisibility and Associates

For the rest of this chapter, we will focus attention on polynomial rings $F[x]$ where the ring F of coefficients is a *field*. Given $f, g \in F[x]$, we say that g divides f (written $g|f$) iff there exists $q \in F[x]$ with $f = gq$. We also say that g is a *divisor* of f , and that f is a *multiple* of g in $F[x]$. For $f = 0$, taking $q = 0$ shows that $g|0$ for every $g \in F[x]$. If f is nonzero and $g|f$, applying the degree addition formula to the equation $f = gq$ shows that g is nonzero and $\deg(g) \leq \deg(f)$. For nonzero g , g divides f iff the remainder when we divide f by g is zero. This observation provides an algorithmic test for deciding whether or not one polynomial divides another.

Observe that $f|f$ for all $f \in F[x]$, since $f = f \cdot 1$. Suppose $f|g$ and $g|h$ in $F[x]$. Then $f|h$ follows, since $g = fq$ and $h = gs$ (for some $q, s \in F[x]$) implies $h = f(qs)$ with $qs \in F[x]$. Thus, divisibility is a reflexive and transitive relation on $F[x]$. Suppose $f|g$ and $g|f$; does it follow that $f = g$? In general, the answer is no. For, the hypotheses $f|g$ and $g|f$ ensure that $g = fs$ and $f = gt$ for some $s, t \in F[x]$, and therefore $f1 = f = gt = f(st)$. Assuming f is nonzero, we can cancel f in the integral domain $F[x]$ to conclude that $st = 1$. By calculating degrees, we see that s and t are nonzero elements of the field F , and $t = s^{-1}$. Thus, $f = gt$ is a scalar multiple of g in the vector space $F[x]$, but f need not *equal* g . On the other hand, if f and g are both known to be *monic*, then $f|g$ and $g|f$ does imply $f = g$, since comparison of leading terms in the equation $f = gt$ forces $t = 1$.

In general, we can always multiply a polynomial in $F[x]$ by a nonzero constant from F (i.e., a polynomial of degree zero) without affecting divisibility properties. More precisely, suppose $f, g \in F[x]$ and $c \in F$ where $c \neq 0$. Then one checks that $f|g$ iff $(cf)|g$ iff $f|(cg)$. For any $f, h \in F[x]$, let us write $f \sim h$ iff there exists a nonzero $c \in F$ with $h = cf$. In this situation, we say that f and h are *associates* in the ring $F[x]$. One verifies that the relation \sim is reflexive, symmetric, and transitive (symmetry requires the hypothesis that F is a field), so \sim is an equivalence relation. Furthermore, the equivalence class of any nonzero f (which is the set of associates of f in $F[x]$) contains exactly one monic polynomial, which can be found by dividing f by its leading coefficient.

Suppose $f, p, q, a, b \in F[x]$ are polynomials such that $f|p$ and $f|q$; then $f|(ap + bq)$. To prove this, write $p = fc$ and $q = fd$ for some $c, d \in F[x]$, and then compute

$$ap + bq = a(fc) + b(fd) = f(ac + bd),$$

where $ac + bd \in F[x]$. It follows by induction that for all $n \geq 1$ and all $f, p_1, \dots, p_n, a_1, \dots, a_n \in F[x]$, if f divides every p_i , then f also divides the “polynomial combination” $a_1p_1 + \dots + a_np_n$.

3.8 Greatest Common Divisors of Polynomials

Suppose f and g are polynomials in $F[x]$ that are not both zero. A *common divisor* of f and g is a polynomial $p \in F[x]$ such that $p|f$ and $p|g$. A *greatest common divisor* (gcd) of f and g is a common divisor p of f and g whose degree is as large as possible. We claim that gcds of f and g always exist, although they may not be unique. To see this, let $\text{CDiv}(f, g)$ denote the set of all common divisors of f and g in $F[x]$. This set is nonempty, since $1|f$ and $1|g$, and the set does not contain zero, since $f \neq 0$ or $g \neq 0$. If $f \neq 0$, then every common divisor has degree at most $\deg(f)$, and similarly if $g \neq 0$. Thus, the set of degrees of polynomials in $\text{CDiv}(f, g)$ is a finite nonempty subset of \mathbb{N} , which must have a greatest element. Any polynomial in $\text{CDiv}(f, g)$ having this maximum degree will be a gcd of f and g .

We now prove the following more precise result: *if F is a field and $f, g \in F[x]$ are not both zero, then there exists a unique monic gcd of f and g , which is denoted $\text{gcd}(f, g)$. The set of all gcds of f and g consists of the associates of $\text{gcd}(f, g)$. Every common divisor of f and g divides $\text{gcd}(f, g)$ in $F[x]$, and conversely. There exist polynomials $u, v \in F[x]$ with $\text{gcd}(f, g) = uf + vg$.* The proof we give, called *Euclid’s Algorithm*, provides a specific computational procedure for finding $\text{gcd}(f, g)$ and polynomials u, v such that $\text{gcd}(f, g) = uf + vg$.

The key to the proof is the following observation. Suppose $a, b, q, r \in F[x]$, $b \neq 0$, and $a = bq + r$; then $\text{CDiv}(a, b) = \text{CDiv}(b, r)$. To see this, assume $p \in F[x]$ satisfies $p|a$ and $p|b$; then p divides $1a + (-q)b = r$. So $\text{CDiv}(a, b) \subseteq \text{CDiv}(b, r)$. Conversely, if $p|b$ and $p|r$, then p divides $qb + 1r = a$. It follows that a polynomial $h \in F[x]$ is a gcd of a and b iff h is a gcd of b and r .

We now describe the algorithm for computing a gcd of f and g . Switching f and g if needed, we can assume that $g \neq 0$. We will construct two finite sequences of polynomials $r_0, r_1, r_2, \dots, r_{N+1}$ and q_1, q_2, \dots, q_N by repeatedly performing polynomial long division. To start, let $r_0 = f$ and $r_1 = g \neq 0$. Divide r_0 by r_1 to obtain

$$r_0 = q_1 r_1 + r_2, \quad (r_2 = 0 \text{ or } \deg(r_2) < \deg(r_1)).$$

If $r_2 = 0$, the construction ends with $N = 1$. If $r_2 \neq 0$, we continue by dividing r_1 by r_2 , obtaining

$$r_1 = q_2 r_2 + r_3, \quad (r_3 = 0 \text{ or } \deg(r_3) < \deg(r_2)).$$

If $r_3 = 0$, the construction ends with $N = 2$. If $r_3 \neq 0$, we continue similarly. At the i 'th step, assuming that r_i is nonzero, we divide r_{i-1} by r_i to obtain

$$r_{i-1} = q_i r_i + r_{i+1}, \quad (r_{i+1} = 0 \text{ or } \deg(r_{i+1}) < \deg(r_i)).$$

If $r_{i+1} = 0$, we end the construction with $N = i$. Otherwise, we proceed to step $i + 1$ of the computation.

We claim first that the computation must end after a finite number of steps. We prove this claim by a contradiction argument. If the algorithm never terminates, then r_i is never zero for any i , so we obtain an infinite decreasing sequence of degrees

$$\deg(r_1) > \deg(r_2) > \deg(r_3) > \dots > \deg(r_i) > \dots . \tag{3.1}$$

Then the set $S = \{\deg(r_i) : i \geq 1\}$ has no least element. But S is a nonempty subset of the well-ordered set \mathbb{N} , so S must have a least element. This contradiction shows that the algorithm does in fact terminate. Indeed, one can prove using (3.1) that $N \leq \deg(r_1) + 1$.

We now analyze the common divisors of f and g . By the observation preceding the description of the algorithm, the relation $r_{i-1} = q_i r_i + r_{i+1}$ implies that $\text{CDiv}(r_{i-1}, r_i) = \text{CDiv}(r_i, r_{i+1})$ for $1 \leq i \leq N$. Accordingly,

$$\text{CDiv}(f, g) = \text{CDiv}(r_0, r_1) = \text{CDiv}(r_1, r_2) = \cdots = \text{CDiv}(r_N, r_{N+1}) = \text{CDiv}(r_N, 0).$$

Thus, the common divisors of f and g are the same as the common divisors of r_N and zero. Since everything divides zero, these are the same as the divisors of r_N . Now, the divisors of r_N of maximum degree are r_N and its associates. Therefore, a polynomial h is a gcd of f and g iff h is an associate of r_N . As pointed out earlier, r_N has a unique *monic* associate, which is denoted $\gcd(f, g)$. Our argument has shown that every common divisor of f and g divides r_N , hence divides every associate of r_N . In particular, all common divisors of f and g divide $\gcd(f, g)$. Conversely, any divisor of $\gcd(f, g)$ is a common divisor of f and g .

To finish the proof, we must find polynomials $u, v \in F[x]$ such that $\gcd(f, g) = uf + vg$. It suffices to construct polynomials $u_i, v_i \in F[x]$ such that $r_i = u_i f + v_i g$ for $0 \leq i \leq N$; for u and v can then be found by dividing u_N and v_N by the leading coefficient of r_N . We argue by strong induction on i . If $i = 0$, choose $u_i = 1$ and $v_i = 0$ (recalling that $r_0 = f$). If $i = 1$, choose $u_i = 0$ and $v_i = 1$ (recalling that $r_1 = g$). For the induction step, assume $1 \leq i \leq N$ and the polynomials $u_{i-1}, v_{i-1}, u_i, v_i$ with the required properties have already been computed. Then

$$\begin{aligned} r_{i+1} &= r_{i-1} - q_i r_i \\ &= (u_{i-1}f + v_{i-1}g) - q_i(u_i f + v_i g) \\ &= (u_{i-1} - q_i u_i)f + (v_{i-1} - q_i v_i)g. \end{aligned}$$

Therefore, $r_{i+1} = u_{i+1}f + v_{i+1}g$ will hold if we choose

$$u_{i+1} = u_{i-1} - q_i u_i \in F[x] \text{ and } v_{i+1} = v_{i-1} - q_i v_i \in F[x].$$

These equations provide an explicit recursive algorithm for computing u and v from f and g .

We do not assert that the polynomials u and v such that $\gcd(f, g) = uf + vg$ are unique. Indeed, for any polynomial $z \in F[x]$, $uf + vg = (u + gz)f + (v - fz)g$, so that there are many possible choices for u and v .

3.9 GCDs of Lists of Polynomials

We now extend the discussion of gcds to the case of more than two polynomials. Suppose $n \geq 1$ and f_1, \dots, f_n are polynomials in $F[x]$. There is no real loss of generality in assuming that all f_i 's are nonzero. We call $g \in F[x]$ a *common divisor* of the f_i 's if g divides every f_i ; let $\text{CDiv}(f_1, \dots, f_n)$ be the set of these polynomials. A *greatest common divisor* of f_1, \dots, f_n is a polynomial of maximum degree in $\text{CDiv}(f_1, \dots, f_n)$. As in the case $n = 2$, one checks that greatest common divisors of the f_i 's exist, although they need not be unique. In the case $n = 1$, $\text{CDiv}(f_1)$ is the set of all divisors of f_1 , the greatest common divisors of the collection $\{f_1\}$ are the associates of f_1 , and $\gcd(f_1)$ is the unique monic associate of f_1 .

We will prove that *for all $n \geq 1$ and all nonzero $f_1, \dots, f_n \in F[x]$, there exists*

a unique monic greatest common divisor of f_1, \dots, f_n , denoted $\gcd(f_1, \dots, f_n)$. Also, $\text{CDiv}(f_1, \dots, f_n) = \text{CDiv}(\gcd(f_1, \dots, f_n))$ (this says that the common divisors of the f_i 's are precisely the divisors of the gcd of the f_i 's); there exist polynomials $u_1, \dots, u_n \in F[x]$ such that $\gcd(f_1, \dots, f_n) = u_1 f_1 + \dots + u_n f_n$; and for all $n \geq 2$,

$$\gcd(f_1, \dots, f_{n-1}, f_n) = \gcd(\gcd(f_1, \dots, f_{n-1}), f_n). \quad (3.2)$$

These assertions hold when $n = 1$, and we have already proved them when $n = 2$.

Proceeding by induction, assume $n > 2$ and the results are already known for collections of $n - 1$ polynomials. We first prove (3.2) as follows. For a fixed $h \in F[x]$, the following conditions are logically equivalent: (a) h divides f_1, \dots, f_{n-1}, f_n ; (b) h divides f_1, \dots, f_{n-1} and $h|f_n$; (c) h divides $\gcd(f_1, \dots, f_{n-1})$ and $h|f_n$; (d) $h|\gcd(\gcd(f_1, \dots, f_{n-1}), f_n)$. The equivalence of (b) and (c) follows from part of the induction hypothesis for lists of $n - 1$ polynomials, and the equivalence of (c) and (d) follows similarly using the theorem for a list of two polynomials. From the equivalence of (a) and (d), we see that $\text{CDiv}(f_1, \dots, f_n)$ is exactly the set of polynomials dividing $g = \gcd(\gcd(f_1, \dots, f_{n-1}), f_n)$. It now follows from the definition that the greatest common divisors of f_1, \dots, f_n consist of g and its associates. In particular, g is the unique monic gcd of f_1, \dots, f_n .

To prove the remaining statement in the theorem, use the induction hypothesis to write

$$\gcd(f_1, \dots, f_{n-1}) = v_1 f_1 + \dots + v_{n-1} f_{n-1}$$

for some polynomials $v_i \in F[x]$. Putting $g = \gcd(\gcd(f_1, \dots, f_{n-1}), f_n)$ as above, we can also write

$$g = c \gcd(f_1, \dots, f_{n-1}) + d f_n$$

for some polynomials $c, d \in F[x]$. Combining these expressions, we see that

$$\gcd(f_1, \dots, f_n) = g = (cv_1)f_1 + \dots + (cv_{n-1})f_{n-1} + df_n.$$

Thus, the gcd is expressible as a polynomial combination of the f_i 's. With the help of Euclid's algorithm for computing the gcd of two polynomials, one may convert this proof into an algorithm that computes $\gcd(f_1, \dots, f_n)$ and the polynomials u_1, \dots, u_n satisfying $\gcd(f_1, \dots, f_n) = \sum_{i=1}^n u_i f_i$. An alternative algorithm to compute this information is presented in the next section.

3.10 Matrix Reduction Algorithm for GCDs

Given polynomials $f, g \in F[x]$, one can compute $d = \gcd(f, g)$ and $a, b \in F[x]$ with $d = af + bg$ using a matrix reduction technique similar to the Gaussian elimination algorithm learned in elementary linear algebra. One begins with the 2×3 matrix $\left[\begin{array}{cc|c} 1 & 0 & f \\ 0 & 1 & g \end{array} \right]$. One then performs a sequence of “elementary row operations” (described below) to convert this matrix into a new matrix of the form $\left[\begin{array}{cc|c} a & b & d \\ r & s & 0 \end{array} \right]$, where all entries are elements of $F[x]$ and one entry in the rightmost column is 0. We will see that d is a gcd of f and g , and $d = af + bg$.

The allowable row operations are as follows. First, one can multiply a row by a nonzero constant in F . Second, one can switch the two rows. Third, for any polynomial $z \in F[x]$, one can add z times row 1 to row 2, or add z times row 2 to row 1.

We illustrate the algorithm by computing $\gcd(x^2 + 2x - 3, x^3 - 2x^2 + 2x - 1)$ in $\mathbb{Q}[x]$. We begin by adding $-x$ times row 1 to row 2 to eliminate the x^3 term:

$$\left[\begin{array}{cc|cc} 1 & 0 & x^2 + 2x - 3 \\ 0 & 1 & x^3 - 2x^2 + 2x - 1 \end{array} \right] \rightarrow \left[\begin{array}{cc|cc} 1 & 0 & x^2 + 2x - 3 \\ -x & 1 & -4x^2 + 5x - 1 \end{array} \right].$$

Next we add 4 times row 1 to row 2:

$$\left[\begin{array}{cc|cc} 1 & 0 & x^2 + 2x - 3 \\ -x & 1 & -4x^2 + 5x - 1 \end{array} \right] \rightarrow \left[\begin{array}{cc|cc} 1 & 0 & x^2 + 2x - 3 \\ -x + 4 & 1 & 13x - 13 \end{array} \right].$$

Now multiply row 2 by $1/13$, then add $-x$ times row 2 to row 1:

$$\rightarrow \left[\begin{array}{cc|cc} 1 & 0 & x^2 + 2x - 3 \\ (-x + 4)/13 & 1/13 & x - 1 \end{array} \right] \rightarrow \left[\begin{array}{cc|cc} (1/13)x^2 - (4/13)x + 1 & -x/13 & 3x - 3 \\ (-x + 4)/13 & 1/13 & x - 1 \end{array} \right].$$

Finally, add -3 times row 2 to row 1 to obtain the matrix

$$\left[\begin{array}{cc|cc} (1/13)x^2 - (1/13)x + (1/13) & (-x - 3)/13 & 0 \\ (-x + 4)/13 & 1/13 & x - 1 \end{array} \right].$$

Switching rows 1 and 2, we conclude that the gcd is $x - 1$, and (as one readily checks)

$$((-x + 4)/13) \cdot (x^2 + 2x - 3) + (1/13) \cdot (x^3 - 2x^2 + 2x - 1) = x - 1.$$

Why does the matrix reduction algorithm work correctly? First, we show that the algorithm will always produce a zero in the rightmost column in finitely many steps if we use appropriate row operations. Suppose the current matrix looks like $\left[\begin{array}{cc|c} a & b & c \\ r & s & t \end{array} \right]$. If c and t are both nonzero, we can argue by induction on $\deg(c) + \deg(t)$. As in the proof of the polynomial division theorem, if $\deg(c) \leq \deg(t)$, we can add an appropriate polynomial multiple of row 1 to row 2 to lower the degree of t . Similarly, if $\deg(c) > \deg(t)$, we can add a multiple of row 2 to row 1 to lower the degree of c . In either case, the nonnegative integer $\deg(c) + \deg(t)$ strictly decreases as a result of the row operation (or one of c or t becomes zero, causing termination). Since there is no infinite strictly decreasing sequence of nonnegative integers, the algorithm will terminate in finite time.

To see why the algorithm gives correct results, we prove the following statement by induction on the number of row operations performed. If the current matrix is $A = \left[\begin{array}{cc|c} a & b & c \\ r & s & t \end{array} \right]$, then $af + bg = c$, $rf + sg = t$, and there exist $u, v, w, y \in F[x]$ with $f = uc + vt$ and $g = wc + yt$. When the algorithm starts, we have $a = s = 1$, $b = r = 0$, $c = f$, and $t = g$. So the given statement holds, taking $u = y = 1$ and $v = w = 0$. For the induction step, assume the statement holds for the matrix A shown above. Suppose we multiply row 1 by a nonzero constant $k \in F$, producing the new matrix $B = \left[\begin{array}{cc|c} ka & kb & kc \\ r & s & t \end{array} \right]$. Then $rf + sg = t$ still holds; $(ka)f + (kb)g = (kc)$ follows from $af + bg = c$; and $f = (uk^{-1})(kc) + vt$, $g = (wk^{-1})(kc) + yt$. We argue similarly if row 2 is multiplied by k . Now suppose we modify the matrix A by adding z times row 2 to row 1, for some $z \in F[x]$. The new matrix is $C = \left[\begin{array}{cc|c} a + zr & b + zs & c + zt \\ r & s & t \end{array} \right]$. Adding z times $rf + sg = t$ to $af + bg = c$, we deduce $(a + zr)f + (b + zs)g = c + zt$, and the equation $rf + sg = t$ is still true. Furthermore, given that $f = uc + vt$ and $g = wc + yt$, we have $f = u(c + zt) + (v - uz)t$ and $g = w(c + zt) + (y - wz)t$ where $u, v - uz, w, y - wz \in F[x]$. Thus, the induction hypothesis holds for the new matrix

C. We argue similarly in the case where we add z times row 1 to row 2. Finally, one may check that the induction hypothesis is preserved if we interchange the two rows of A .

Now consider what happens at the end of the algorithm, when one of the entries in the rightmost column of A (say t) becomes zero. Since $af + bg = c$, we see that any common divisor of f and g must also divide c . On the other hand, since $f = uc + v0 = uc$ and $g = wc + y0 = wc$, c is a common divisor of f and g . Thus, c must be a gcd of f and g . We can use one more row operation to make c monic.

We can generalize the matrix reduction algorithm to take given nonzero inputs $f_1, \dots, f_n \in F[x]$ (for any fixed $n \geq 1$) and compute $d = \gcd(f_1, \dots, f_n)$ and $u_1, \dots, u_n \in F[x]$ with $d = u_1f_1 + \dots + u_nf_n$. One starts with an $n \times n$ identity matrix, augmented by an extra column with entries f_1, \dots, f_n . Performing elementary row operations, one reduces this matrix to a new one with only one nonzero entry in the last column. This entry is a gcd of the f_i 's, and the remaining entries in its row give the coefficients u_i expressing d as a polynomial combination of the f_i 's. We ask the reader to provide a detailed proof of this assertion (Exercise 41), which is similar to the proof given above for $n = 2$.

The matrix reduction algorithm can also be adapted to the computation of gcds of two or more integers. In this case, one is only allowed to multiply a particular row by 1 or -1 , but one can add any integer multiple of a row to a different row. To prove termination, one can use induction on the sum of the absolute values of the entries in the rightmost column. Examples and proofs for the integer version of this algorithm appear in the exercises.

3.11 Roots of Polynomials

Let F be a field. Given $p \in F[x]$ and $c \in F$, we call c a *root* or *zero* of p iff $p(c) = f_p(c) = 0$. We now show that c is a root of p iff $(x - c)$ divides p in the ring $F[x]$. To see this, fix $p \in F[x]$ and $c \in F$. Dividing p by $x - c$ gives $p = q(x - c) + r$ for some $q, r \in F[x]$ with $r = 0$ or $\deg(r) < \deg(x - c) = 1$. This means that the remainder r is a constant polynomial (possibly zero). Evaluating both sides of $p = q(x - c) + r$ at $x = c$ shows that $p(c) = q(c)(c - c) + r = r$. Now, c is a root of p iff $p(c) = 0$ iff the remainder $r = 0$ iff $(x - c)|p$ in $F[x]$.

More generally, if K is a field with subfield F , $p \in F[x]$, and $c \in K$, then c is called a *root* of p iff $p(c) = 0$. If $c \notin F$, the previous theorem does not apply directly, since $x - c$ is not an element of $F[x]$. However, since $F \subseteq K$, we can regard p as an element of the polynomial ring $K[x]$. Then c is a root of p iff $(x - c)|p$ in the ring $K[x]$.

We can now establish a bound on the number of roots of a polynomial. *If F is a field and $p \in F[x]$ has degree $n > 0$, then p has at most n roots in F .* We prove this result by induction on n . If $n = 1$, write $p = ax + b$ for some $a, b \in F$ with $a \neq 0$. One may check that $x = -a^{-1}b$ is the unique root of this first-degree polynomial. Now suppose $n > 1$ and the result is known for polynomials of degree $n - 1$. Given p of degree n , consider two cases. If p has no root in F , then the result certainly holds. Otherwise, let c be a root of p in F . Then we can write $p = (x - c)q$ for some polynomial q of degree $n - 1$. If $d \in F$ is any root of p different from c , then $0 = p(d) = (d - c)q(d)$. Since $d - c \neq 0$ and F is a field, we see that $q(d) = 0$. This means that every root of p besides c is a root of q ; conversely, every root of q is also a root of p . Now, by the induction hypothesis, q has at most $n - 1$ roots in F (possibly including c). Thus, p has at most n roots in F , namely all the roots of q together with c .

The preceding theorem would still work if we allowed polynomials $p \in R[x]$ where R is

an integral domain. But the theorem can fail dramatically for general commutative rings R . For instance, the polynomial $p = x^2 + 7 \in \mathbb{Z}_8[x]$ has four roots (namely 1, 3, 5, and 7) in the ring $R = \mathbb{Z}_8$. Note that this ring is neither an integral domain nor a field.

Next let F be any field, and consider the question of finding the roots of a given $p \in F[x]$ in the field F . If F is a finite field, we can find the roots of p (if any) by evaluating p at all the elements of F and seeing when the answer is zero. Now let $F = \mathbb{Q}$ be the field of rational numbers. Given a non-constant $p \in \mathbb{Q}[x]$, we can replace p by an associate cp (for some nonzero $c \in \mathbb{Q}$) such that cp is in $\mathbb{Z}[x]$, and cp has the same rational roots as p . To accomplish this, we could pick c to be the product (or least common multiple) of all denominators appearing in coefficients of p . We can now write $cp = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$ for certain integers a_0, a_1, \dots, a_n with $a_n \neq 0$. We will prove the *rational root theorem*, which asserts that *every rational root of p must have the form r/s for some integers r, s such that $r|a_0$ and $s|a_n$* .

To begin the proof, assume r, s are relatively prime integers such that r/s is a root of p (hence also a root of cp). Evaluating cp at this root gives the equation

$$0 = a_n(r^n/s^n) + a_{n-1}(r^{n-1}/s^{n-1}) + \cdots + a_1(r^1/s^1) + a_0.$$

Now multiply through by the integer s^n to clear the fractions, giving

$$0 = a_n r^n + a_{n-1} r^{n-1} s^1 + \cdots + a_1 r^1 s^{n-1} + a_0 s^n.$$

Isolating a_0 on one side, we obtain

$$a_0 s^n = -(a_n r^n + a_{n-1} r^{n-1} s^1 + \cdots + a_1 r^1 s^{n-1}).$$

The right side is a multiple of r in \mathbb{Z} , so $r|a_0 s^n$ in \mathbb{Z} . As r and s are relatively prime, $r|a_0$ follows. (This can be shown, for instance, by considering prime factorizations; cf. Exercise 33.) Similarly, isolating a_n in the earlier equation shows that $a_n r^n$ is a multiple of s , so $s|a_n r^n$, so $s|a_n$ because $\gcd(r, s) = 1$.

To illustrate the rational root theorem, consider the polynomial $p = x^4 + (13/3)x^3 + (22/3)x^2 + 6x + (4/3)$ in $\mathbb{Q}[x]$. Taking $c = 3$, we have $cp = 3x^4 + 13x^3 + 22x^2 + 18x + 4 \in \mathbb{Z}[x]$. We find $a_4 = 3$ and $a_0 = 4$, so the possible rational roots of p are $\pm 1, \pm 1/3, \pm 2, \pm 2/3, \pm 4$, and $\pm 4/3$. Evaluating p at each of these points, we find that p has rational roots $-1/3$ and -2 . We can then factor $p = (x + 1/3)(x + 2)(x^2 + 2x + 2)$ in $\mathbb{Q}[x]$. As another example, let $q = x^2 - 3 \in \mathbb{Q}[x]$. By the rational root theorem, the only possible roots of q in \mathbb{Q} are ± 1 and ± 3 . Since none of these four numbers is a root of q , we can deduce that $\sqrt{3}$ is not rational. For a generalization of this example, see Exercise 64.

3.12 Irreducible Polynomials

Prime numbers and factorizations of numbers into products of primes play a pivotal role in the theory of integers. Our next goal is to develop analogous concepts for the theory of polynomials.

Let F be a field, and let $p \in F[x]$ have positive degree. We shall call p an *irreducible polynomial over F* iff whenever $p = fg$ for some polynomials f and g in $F[x]$, either f or g has degree zero. In this case, if p has degree $n > 0$ and f (say) has degree zero, then g has degree n and is an associate of p . Conversely, for any nonzero scalar $c \in F$, we have the factorization $p = c(p/c)$ where p/c is an associate of p . It follows that the divisors of an

irreducible polynomial p are precisely the nonzero constant polynomials and the associates of p . One can check that an associate of p is irreducible iff p is irreducible.

Negating the definition of irreducibility, we see that p is *reducible* in $F[x]$ iff there exists a “proper” factorization $p = fg$, where neither f nor g is constant. Computing degrees, we see that we must have $0 < \deg(f) < \deg(p)$ and $0 < \deg(g) < \deg(p)$ in this situation. This leads to a convenient criterion for detecting irreducibility of polynomials p of degree at most 3. First, every polynomial of degree 1 is irreducible, since otherwise $\deg(f)$ would be an integer strictly between 0 and 1. Second, if $p \in F[x]$ is such that $\deg(p) > 1$ and $p(c) = 0$ for some $c \in F$, then $(x - c)|p$ in $F[x]$ and hence p is reducible. Conversely, if $\deg(p)$ is 2 or 3 and p is reducible, then one of the factors in a proper factorization $p = fg$ must have degree 1. Moving a constant from f to g if needed, we can assume p has a *monic* factor $x - c$, and then $c \in F$ is a root of p . Thus, *a polynomial p of degree 2 or 3 is irreducible in $F[x]$ if and only if p has no root in F . Irreducible polynomials of larger degree have no root in F , but the converse does not hold in general.* For example, $p = (x^2 - 3)^2$ is visibly reducible in $\mathbb{Q}[x]$, yet p has no root in \mathbb{Q} .

Note that the field of coefficients F must be considered when determining irreducibility, since we only allow factorizations $p = fg$ where the factors f and g lie in $F[x]$. For example, consider the polynomial $p = x^2 - 2$. Applying the criterion in the last paragraph, we see that p is irreducible in $\mathbb{Q}[x]$ (since there is no *rational* c with $c^2 = 2$), but p is reducible in $\mathbb{R}[x]$, with factorization $x^2 - 2 = (x - \sqrt{2})(x + \sqrt{2})$. Similarly, $x^2 + 1$ is irreducible in $\mathbb{Q}[x]$ and in $\mathbb{R}[x]$, but becomes reducible in the ring $\mathbb{C}[x]$, since $i^2 + 1 = 0$.

For some fields F , one can give specific algorithms for testing the irreducibility of polynomials $p \in F[x]$. For example, if F is any *finite* field, such as \mathbb{Z}_q for prime q , then a given $p \in F[x]$ has only a finite number of possible divisors. We can use polynomial long division to test each non-constant $f \in F[x]$ of degree less than $n = \deg(p)$ to see if $f|p$; p is irreducible in $F[x]$ iff no such divisor f exists. In fact, one can use the degree addition formula to verify that if $p = fg$ is reducible, then one of its factors must have degree at most $\lfloor n/2 \rfloor$, so we need only test divisors up through this degree. The same process gives a factoring algorithm to find all possible divisors of a reducible $f \in F[x]$.

Irreducibility over infinite fields is usually harder to check, with a few exceptions. For the field \mathbb{C} , the famous *fundamental theorem of algebra* asserts that every non-constant $p \in \mathbb{C}[x]$ has a root in \mathbb{C} . We will not prove this theorem here. But combining this theorem with our earlier observations relating roots to irreducibility, we see that *$p \in \mathbb{C}[x]$ is irreducible in $\mathbb{C}[x]$ iff $\deg(p) = 1$.* With a little more work, one can also use the fundamental theorem of algebra to deduce that *$p \in \mathbb{R}[x]$ is irreducible in $\mathbb{R}[x]$ iff $\deg(p) = 1$ or $p = ax^2 + bx + c$ for some $a, b, c \in \mathbb{R}$ with $a \neq 0$ and $b^2 - 4ac < 0$.* Later in this chapter, we will prove some theorems and give an algorithm that can be used to test irreducibility of polynomials in $\mathbb{Q}[x]$.

Let us return to the case of a general field F . We now establish a crucial property of irreducible polynomials, which is an analogue of a corresponding property of prime integers: *for all $p, r, s \in F[x]$, if p is irreducible and $p|rs$, then $p|r$ or $p|s$.* To prove this, fix $p, r, s \in F[x]$ with p irreducible, and assume $p|rs$ but p does not divide r . Since $p|rs$, we may write $rs = pt$ for some $t \in F[x]$. What is $\gcd(p, r)$? This gcd must be a divisor of the irreducible polynomial p , so it is either 1 or the unique monic associate of p . The latter possibility is ruled out since p does not divide r , so $\gcd(p, r) = 1$. We may therefore write $1 = ap + br$ for some $a, b \in F[x]$. Multiplying by s gives $s = (ap + br)s = aps + brs = aps + bpt = p(as + bt)$. We see from this that $p|s$. More generally, an induction argument using the result just proved shows that *if an irreducible polynomial $p \in F[x]$ divides a product $f_1 f_2 \cdots f_n$, then p divides f_i for some $i \in [n]$.* Conversely, if $p \in F[x]$ is any non-constant polynomial such that for all $r, s \in F[x]$, $p|rs$ implies $p|r$ or $p|s$, then p must be irreducible in $F[x]$. For if p

were reducible, with some proper factorization $p = fg$ ($f, g \in F[x]$), then $p|fg$ (as p divides itself), but p cannot divide f or g since these polynomials have degree less than p .

3.13 Factorization of Polynomials into Irreducibles

We are ready to prove the following fundamental result on factorizations of polynomials in $F[x]$, where F is a field. *Every non-constant polynomial $g \in F[x]$ can be written in the form $g = cp_1p_2 \cdots p_r$, where $c \in F$ is a nonzero constant and each p_i is a monic, irreducible polynomial in $F[x]$. If also $g = dq_1q_2 \cdots q_s$, where $d \in F$ is constant and the q_j 's are monic, irreducible polynomials in $F[x]$, then $c = d$, $r = s$, and (after appropriate reordering) $p_i = q_i$ for $1 \leq i \leq r$.*

To prove the existence of the asserted factorization for g , we use strong induction on $\deg(g)$. Suppose $n = \deg(g) \geq 1$ and the result is already known for all non-constant polynomials of degree less than n . If g happens to be irreducible (which always occurs when $n = 1$), we have the factorization $g = c(g/c)$ where $c \in F$ is the leading coefficient of g . If g is reducible, let $g = fh$ where $f, h \in F[x]$ satisfy $0 < \deg(f) < \deg(g)$ and $0 < \deg(h) < \deg(g)$. The induction hypothesis allows us to write f and h as constants times products of monic irreducible polynomials. Multiplying these expressions, we obtain a similar expression for g by combining the two constants. This completes the existence proof.

To prove uniqueness, we also use strong induction on the degree of g . Suppose $g = cp_1p_2 \cdots p_r = dq_1q_2 \cdots q_s$ as in the theorem statement, and assume the uniqueness result is already known to hold for factorizations of polynomials of degree less than $\deg(g)$. Since the p_i 's and q_j 's are all monic, comparison of leading coefficients immediately gives $c = d$. Multiplying by c^{-1} gives $p_1p_2 \cdots p_r = q_1q_2 \cdots q_s$. We deduce from this equality that p_r divides the product $q_1q_2 \cdots q_s$. Since p_r is irreducible, p_r must divide some q_j . Reordering the q 's, we may assume that p_r divides q_s . Since q_s is irreducible and p_r is non-constant, p_r must be an associate of q_s . But p_r and q_s are both monic, so $p_r = q_s$ is forced. Since $F[x]$ is an integral domain, we can cancel this nonzero factor from $p_1p_2 \cdots p_r = q_1q_2 \cdots q_s$ to get $p_1p_2 \cdots p_{r-1} = q_1q_2 \cdots q_{s-1}$. If $r = 1$, then the empty product on the left side is 1, which forces $s = 1 = r$. Similarly, $s = 1$ forces $r = 1 = s$. In the remaining case, we have a nonconstant polynomial $h = p_1p_2 \cdots p_{r-1} = q_1q_2 \cdots q_{s-1}$. Since $\deg(h) < \deg(g)$, the induction hypothesis is applicable to these two factorizations of h . We conclude that $r - 1 = s - 1$ and (after reordering) $p_i = q_i$ for $1 \leq i \leq r - 1$. Therefore, $r = s$, $p_i = q_i$ for $1 \leq i \leq r - 1$, and $p_r = q_s = q_r$. This completes the induction.

We can rephrase the preceding result in the following way. Given a field F , let $\{p_i : i \in I\}$ be the collection of all monic irreducible polynomials in $F[x]$, where I is some indexing set. Every non-constant $f \in F[x]$ can be written uniquely in the form

$$f = c \prod_{i \in I} p_i^{e_i}, \quad (3.3)$$

where c is a nonzero element of F , the e_i 's are nonnegative integers, and all but a finite number of e_i 's are zero. This expression for f is obtained from the previous theorem by writing $f = cq_1q_2 \cdots q_s$ (where the q_j 's are monic and irreducible in $F[x]$) and then letting e_i be the number of q_j 's equal to p_i . Uniqueness of the e_i 's follows immediately from the uniqueness of the factorization of f (note that reordering the q_j 's will not affect the exponents e_i). If we allow *all* e_i 's to be zero, we see that every nonzero polynomial (even a constant) can be written uniquely in the form (3.3). The zero polynomial can also be

written in this form if we allow $c = 0$, but the e_i 's are not unique in this case (typically we would take all e_i 's to be zero). The expression (3.3) is often called the *prime factorization of the polynomial f* . We sometimes write $c = c_f$ and $e_i = e_i(f)$ to indicate the dependence of these quantities on f .

3.14 Prime Factorizations and Divisibility

Now that unique prime factorizations of polynomials are available, we can approach the theory of divisibility and polynomial gcds from a new viewpoint. Let F be a field, and let f, g be nonzero polynomials in $F[x]$. Assume first that $g|f$ in $F[x]$, so that $f = gh$ for some polynomial $h \in F[x]$. Writing the prime factorizations (3.3) of f , g , and h , the relation $f = gh$ becomes

$$c_f \prod_{i \in I} p_i^{e_i(f)} = \left(c_g \prod_{i \in I} p_i^{e_i(g)} \right) \cdot \left(c_h \prod_{i \in I} p_i^{e_i(h)} \right) = (c_g c_h) \prod_{i \in I} p_i^{e_i(g) + e_i(h)},$$

where $c_f, c_g, c_h \in F$ and $e_i(f), e_i(g), e_i(h) \in \mathbb{N}$. By uniqueness of the prime factorization of f , it follows that $c_f = c_g c_h$ and $e_i(f) = e_i(g) + e_i(h)$ for all $i \in I$. Since $e_i(h) \geq 0$, we see that $e_i(g) \leq e_i(f)$ for all $i \in I$. Conversely, if $e_i(g) \leq e_i(f)$ for all $i \in I$, then $g|f$ since

$$f = g \cdot (c_f/c_g) \prod_{i \in I} p_i^{e_i(f) - e_i(g)}.$$

So for all nonzero $f, g \in F[x]$, g divides f iff $e_i(g) \leq e_i(f)$ for all $i \in I$.

Next consider nonzero polynomials $f_1, \dots, f_n, g \in F[x]$. By the preceding paragraph, we see that g is a common divisor of f_1, \dots, f_n iff $e_i(g) \leq e_i(f_j)$ for all $i \in I$ and all $j \in [n]$ iff $e_i(g) \leq \min(e_i(f_1), \dots, e_i(f_n))$ for all $i \in I$. To obtain a common divisor g of the largest possible degree, we need to take $e_i(g) = \min(e_i(f_1), \dots, e_i(f_n))$ for all $i \in I$. This means that

$$\gcd(f_1, \dots, f_n) = \prod_{i \in I} p_i^{\min(e_i(f_1), \dots, e_i(f_n))}.$$

The previous discussion also gives a new proof that the common divisors of the f_j coincide with the divisors of $\gcd(f_1, \dots, f_n)$.

A *common multiple* of $f_1, \dots, f_n \in F[x]$ is a polynomial $h \in F[x]$ such that every f_j divides h . A *least common multiple* of f_1, \dots, f_n is a common multiple of the f_j having minimum degree. Note that h is a common multiple of the f_j iff $e_i(f_j) \leq e_i(h)$ for all $j \in [n]$ and all $i \in I$ iff $\max(e_i(f_1), \dots, e_i(f_n)) \leq e_i(h)$ for all $i \in I$. To keep the degree of h as small as possible, we need to take $e_i(h) = \max(e_i(f_1), \dots, e_i(f_n))$ for all $i \in I$. It is now evident that the least common multiples of f_1, \dots, f_n are precisely the associates of the monic polynomial

$$\text{lcm}(f_1, \dots, f_n) = \prod_{i \in I} p_i^{\max(e_i(f_1), \dots, e_i(f_n))}.$$

Furthermore, the common multiples of f_1, \dots, f_n are precisely the multiples of $\text{lcm}(f_1, \dots, f_n)$. In the case $n = 2$, we have $\max(e_i(f_1), e_i(f_2)) + \min(e_i(f_1), e_i(f_2)) = e_i(f_1) + e_i(f_2)$ for all $i \in I$. Comparing prime factorizations, we deduce the identity

$$\text{lcm}(f_1, f_2) \gcd(f_1, f_2) = c f_1 f_2$$

for some constant $c \in F$.

3.15 Irreducible Polynomials in $\mathbb{Q}[x]$

Let us now return to the question of how to decide whether a specific polynomial $f \in \mathbb{Q}[x]$ is irreducible, where \mathbb{Q} is the field of rational numbers. We have remarked earlier that f is irreducible iff any associate of f is irreducible. By replacing f by one of its associates, we can assume that f lies in $\mathbb{Z}[x]$, i.e., all the coefficients of f are integers. For example, we can accomplish this by multiplying the original f by the product of all the denominators of the coefficients appearing in f . For the rest of our discussion, we will assume f is in $\mathbb{Z}[x]$, and we want to know if f is irreducible in $\mathbb{Q}[x]$.

Recall that all degree 1 polynomials are irreducible. Also, for f of degree 2 or 3, we know f is irreducible in $\mathbb{Q}[x]$ iff f has a rational root. For $f \in \mathbb{Z}[x]$ of arbitrary degree, we can find all rational roots of f (if any) by checking the finitely many rational numbers satisfying the condition in the rational root theorem (§3.11).

To continue, we need the following rather technical lemma. *Assume $f \in \mathbb{Z}[x]$ has degree $n > 0$ and is reducible in $\mathbb{Q}[x]$. Then there exist g, h in $\mathbb{Z}[x]$ with $f = gh$ and $\deg(g), \deg(h) < n$.* Note that the conclusion would follow immediately from the definition of reducible polynomials if we allowed g, h to come from $\mathbb{Q}[x]$, but we are demanding the stronger conclusion that g and h lie in $\mathbb{Z}[x]$. Obtaining this stronger result is surprisingly tricky, so we proceed in several steps.

Step 1: We show that for all prime integers p and all $g, h \in \mathbb{Z}[x]$, if p divides gh in the ring $\mathbb{Z}[x]$, then $p|g$ or $p|h$ in the ring $\mathbb{Z}[x]$. We prove the contrapositive. First, note that $p|g$ in $\mathbb{Z}[x]$ iff $g = pq$ for some $q \in \mathbb{Z}[x]$ iff every coefficient in g is divisible (in \mathbb{Z}) by p . Now, fix a prime $p \in \mathbb{Z}$ and $g, h \in \mathbb{Z}[x]$, and assume p does not divide g and p does not divide h in $\mathbb{Z}[x]$. We must prove that p does not divide gh in $\mathbb{Z}[x]$. Write $g = \sum_{i \geq 0} g_i x^i$, $h = \sum_{j \geq 0} h_j x^j$ for some $g_i, h_j \in \mathbb{Z}$. By our assumption on g and h , at least one g_i and at least one h_j are not divisible (in \mathbb{Z}) by p . Let r be maximal such that p does not divide g_r , and let s be maximal such that p does not divide h_s . Now, observe that the coefficient of x^{r+s} in gh is

$$g_r h_s + \sum_{i,j: i+j=r+s, i>r} g_i h_j + \sum_{i,j: i+j=r+s, j>s} g_i h_j.$$

By maximality of r , each term in the first sum is an integer multiple of p . By maximality of s , each term in the second sum is an integer multiple of p . But, since p is prime and neither g_r nor h_s is a multiple of p , $g_r h_s$ is not an integer multiple of p . We conclude that the coefficient of x^{r+s} in gh is not a multiple of p , and hence p does not divide gh in $\mathbb{Z}[x]$.

Step 2: We show that if $f \in \mathbb{Z}[x]$ has degree $n > 0$ and is reducible in $\mathbb{Q}[x]$, then there exist $g, h \in \mathbb{Z}[x]$ and $d \in \mathbb{N}^+$ with $df = gh$ and $\deg(g), \deg(h) < n$. By reducibility of f in $\mathbb{Q}[x]$, there exist $u, v \in \mathbb{Q}[x]$ with $f = uv$ and $\deg(u), \deg(v) < n$. The idea is to “clear denominators” in u and v . Say $u = \sum_{i=0}^s (a_i/b_i)x^i$ and $v = \sum_{j=0}^t (c_j/d_j)x^j$ for some $s, t \in \mathbb{N}$, $a_i, c_j \in \mathbb{Z}$, and $b_i, d_j \in \mathbb{N}^+$. Choose $d = b_0 b_1 \cdots b_s d_0 d_1 \cdots d_t \in \mathbb{N}^+$, $g = (b_0 \cdots b_s)u$, and $h = (d_0 \cdots d_t)v$. One sees immediately that $g, h \in \mathbb{Z}[x]$, $df = gh$, and $\deg(g), \deg(h) < n$.

Step 3: Among all possible choices of d, g, h satisfying the conclusion of Step 2, choose one with d minimal (which can be done, since \mathbb{N}^+ is well-ordered). We will prove $d = 1$, which will finish the proof of the lemma. To get a contradiction, assume that $d > 1$. Then some prime integer p divides d in \mathbb{Z} , say $d = pd_1$ for some $d_1 \in \mathbb{N}^+$. Now, p divides df in $\mathbb{Z}[x]$, and $df = gh$, so $p|gh$ in $\mathbb{Z}[x]$. By Step 1, $p|g$ or $p|h$ in $\mathbb{Z}[x]$. Say $p|g$ (the other case is similar). Then $g = pg^*$ for some $g^* \in \mathbb{Z}[x]$. Now $df = gh$ becomes $pd_1 f = pg^* h$. As $\mathbb{Z}[x]$ is an integral domain, we can cancel p to obtain $d_1 f = g^* h$. Now d_1, g^* , and h satisfy the conclusion of Step 2. But since $d_1 < d$, this contradicts the minimality of d . The proof is now complete.

3.16 Irreducibility in $\mathbb{Q}[x]$ via Reduction Mod p

Equipped with the lemma from the last section, we are ready to prove two theorems that can be used to prove irreducibility of certain polynomials in $\mathbb{Q}[x]$. The first theorem involves “reducing a polynomial mod p ,” where p is a prime integer. Given $f = \sum_{i \geq 0} f_i x^i \in \mathbb{Z}[x]$ and a prime $p \in \mathbb{Z}$, define $\nu_p(f) = \sum_{i \geq 0} (f_i \bmod p) x^i \in \mathbb{Z}_p[x]$; we call $\nu_p(f)$ the *reduction of f modulo p* . Using the definitions of the algebraic operations on polynomials, one readily confirms that $\nu_p : \mathbb{Z}[x] \rightarrow \mathbb{Z}_p[x]$ is a ring homomorphism; in particular, $\nu_p(gh) = \nu_p(g)\nu_p(h)$ for all $g, h \in \mathbb{Z}[x]$.

Our first irreducibility criterion is as follows. *Suppose $f \in \mathbb{Z}[x]$ has degree $n > 0$ and $p \in \mathbb{Z}$ is a prime not dividing the leading coefficient of f . If $\nu_p(f)$ is irreducible in $\mathbb{Z}_p[x]$, then f is irreducible in $\mathbb{Q}[x]$.* We prove the contrapositive. Fix $f \in \mathbb{Z}[x]$ of degree $n > 0$ and a prime p not dividing f 's leading coefficient, and assume f is reducible in $\mathbb{Q}[x]$. By the lemma of §3.15, we can write $f = gh$ for some $g, h \in \mathbb{Z}[x]$ with $\deg(g), \deg(h) < n$. Reducing $f = gh$ modulo p gives $\nu_p(f) = \nu_p(g)\nu_p(h)$ in $\mathbb{Z}_p[x]$. By hypothesis, the leading coefficient of f does not become zero when we pass to $\nu_p(f)$ by reducing coefficients mod p . So $\nu_p(f)$ still has degree n in $\mathbb{Z}_p[x]$, and $\nu_p(g), \nu_p(h)$ have lower degree than n . So $\nu_p(f)$ is reducible in $\mathbb{Z}_p[x]$.

For example, consider $f = x^4 + 6x^2 + 3x + 7 \in \mathbb{Z}[x]$. Reducing mod 2 gives $\nu_2(f) = x^4 + x + 1 \in \mathbb{Z}_2[x]$. The irreducible polynomials of degree at most two in $\mathbb{Z}_2[x]$ are x , $x + 1$, and $x^2 + x + 1$. None of these divide $x^4 + x + 1$ in $\mathbb{Z}_2[x]$, so the latter polynomial is irreducible in $\mathbb{Z}_2[x]$. As f is monic, our criterion applies to show f is irreducible in $\mathbb{Q}[x]$. Similarly, one may check that $x^5 - 4x + 2$ is irreducible in $\mathbb{Q}[x]$ by verifying that its reduction mod 3, namely $x^5 + 2x + 2$, is irreducible in $\mathbb{Z}_3[x]$. In turn, this follows by exhaustively checking the finitely many possible factors of degree 2 or less. On the other hand, note that $\nu_2(x^5 - 4x + 2) = x^5$ is reducible in $\mathbb{Z}_2[x]$. This shows that for a fixed choice of p , the converse of the irreducibility criterion need not hold. In fact, one can find examples of monic polynomials $f \in \mathbb{Z}[x]$ that are irreducible in $\mathbb{Q}[x]$, and yet $\nu_p(f)$ is reducible in $\mathbb{Z}_p[x]$ for every prime integer p . It can be shown that $x^4 - 10x^2 + 1$ has this property (see [49, Example 26, p. 66]).

3.17 Eisenstein’s Irreducibility Criterion for $\mathbb{Q}[x]$

Our second theorem on irreducibility in $\mathbb{Q}[x]$ is called *Eisenstein’s irreducibility criterion*: given $f = a_0 + a_1 x + \cdots + a_n x^n \in \mathbb{Z}[x]$ with $n > 0$ and $a_n \neq 0$, suppose some prime $p \in \mathbb{Z}$ satisfies $p|a_i$ for all $0 \leq i < n$, $p \nmid a_n$, and $p^2 \nmid a_0$; then f is irreducible in $\mathbb{Q}[x]$.

The proof is by contradiction. With f and p as in the theorem, assume f is reducible in $\mathbb{Q}[x]$. By the lemma of §3.15, there exist $g, h \in \mathbb{Z}[x]$ with $f = gh$ and $\deg(g), \deg(h) < n$. Reducing $f = gh$ modulo p , we get $\nu_p(f) = \nu_p(g)\nu_p(h)$. By hypothesis, $b = a_n \bmod p$ is nonzero, but all other coefficients of f reduce to 0 mod p . So $\nu_p(f) = bx^n \in \mathbb{Z}_p[x]$ with $b \neq 0$. Since p is prime, we know \mathbb{Z}_p is a field, and therefore we can apply our earlier theorem on existence of unique factorization into irreducible factors (§3.13) in the polynomial ring $\mathbb{Z}_p[x]$. As x is evidently irreducible in $\mathbb{Z}_p[x]$, the only way that $\nu_p(f) = bx^n$ can factor into a product $\nu_p(g)\nu_p(h)$ is to have $\nu_p(g) = cx^i$ and $\nu_p(h) = dx^j$ for some $c, d \in \mathbb{Z}_p$ with $cd = b$ and some $i, j \in \mathbb{N}$ with $i + j = n$ and $0 < i, j < n$. But this can only happen if

$$g = c'x^i + \cdots + ps, \quad h = d'x^j + \cdots + pt$$

for some $c', d', s, t \in \mathbb{Z}$, where $c' \bmod p = c$, $d' \bmod p = d$, and all coefficients not shown are integer multiples of p . Then the constant term of $f = gh$ is $a_0 = p^2st$, contradicting the assumption that $p^2 \nmid a_0$.

For example, Eisenstein's criterion applies (with $p = 2$) to prove once again that $x^5 - 4x + 2$ is irreducible in $\mathbb{Q}[x]$. We can see that $x^{10} - 2100$ is irreducible in $\mathbb{Q}[x]$ by applying the criterion with $p = 3$ or $p = 7$, but not with $p = 2$ or $p = 5$. More generally, for all $n \geq 1$ and $a \in \mathbb{Z}$ and prime $p \in \mathbb{Z}$ such that $p|a$ but p^2 does not divide a , the criterion shows that $x^n - a$ is irreducible in $\mathbb{Q}[x]$. With a little extra work (see Exercise 66), Eisenstein's criterion can be used to show that for every prime $p \in \mathbb{Z}$, the polynomial $x^{p-1} + x^{p-2} + \cdots + x^2 + x + 1$ is irreducible in $\mathbb{Q}[x]$.

3.18 Kronecker's Algorithm for Factoring in $\mathbb{Q}[x]$

This section develops an algorithm, due to Kronecker, that takes as input an arbitrary polynomial $f \in \mathbb{Z}[x]$ and returns as output the factorization of f into irreducible factors in $\mathbb{Q}[x]$. In particular, we can use this algorithm to test if a given $f \in \mathbb{Z}[x]$ is irreducible in $\mathbb{Q}[x]$. Before describing this algorithm, we need some background on the Lagrange interpolation formula.

Given any field F and a polynomial $f \in F[x]$ of degree n , we saw in §3.11 that f can have at most n roots in F . Next consider two polynomials $g, h \in F[x]$, each of which is zero or has degree at most n . We assert that if $g(c) = h(c)$ for $n+1$ or more values $c \in F$, then $g = h$ in $F[x]$. To see this, consider $f = g - h \in F[x]$, which is either zero or has degree at most n . The latter possibility cannot occur, since f has at least $n+1$ roots in F . So $f = 0$ and $g = h$.

We can now prove the following theorem that lets us build polynomials that send prescribed inputs to prescribed outputs. *For all fields F , all $n \in \mathbb{N}$, all distinct $a_0, a_1, \dots, a_n \in F$, and all (not necessarily distinct) $b_0, b_1, \dots, b_n \in F$, there exists a unique polynomial $g \in F[x]$ such that $g = 0$ or $\deg(g) \leq n$, and $g(a_i) = b_i$ for $0 \leq i \leq n$.* Uniqueness of g follows from the result in the previous paragraph. To prove existence, first note that for $0 \leq i \leq n$,

$$p_i = \prod_{\substack{j=0 \\ j \neq i}}^n (x - a_j) \Bigg/ \prod_{\substack{j=0 \\ j \neq i}}^n (a_i - a_j)$$

is a polynomial of degree n in $F[x]$ satisfying $p_i(a_i) = 1_F$ and $p_i(a_j) = 0_F$ for all $j \neq i$. (The denominator is a nonzero element of F , hence is invertible, because a_0, \dots, a_n are distinct.) Then $g = \sum_{i=0}^n b_i p_i$ is in $F[x]$, has degree at most n (or is zero), and sends a_i to b_i for all i . The explicit formula for g given here is often called the *Lagrange interpolation formula*.

Now we are ready to describe *Kronecker's algorithm* for factoring in $\mathbb{Q}[x]$. As we have seen, it suffices to consider non-constant input polynomials f lying in $\mathbb{Z}[x]$. Say $\deg(f) = n > 1$. If f is reducible in $\mathbb{Q}[x]$, the lemma of §3.15 shows that $f = gh$ for some $g, h \in \mathbb{Z}[x]$ of smaller degree than n . In fact, the degree addition formula ensures that one of g or h (say g) has degree at most $m = \lfloor n/2 \rfloor$. It will suffice to find g (or, to prove irreducibility of f , to show that g must be a constant). For then we can find h by polynomial division, and apply the algorithm recursively to g and to h until the complete factorization into irreducible polynomials is found.

To proceed, pick arbitrary distinct integers a_0, a_1, \dots, a_m . If $f(a_i) = 0$ for some i , then a_i is a root of f in \mathbb{Z} , and hence $g = x - a_i$ is a divisor of f in $\mathbb{Q}[x]$. Assume henceforth that

$f(a_i) \neq 0$ for every i . The hypothesized factorization $f = gh$ implies $f(a_i) = g(a_i)h(a_i)$ for all i . In these equations, everything is an integer, and each $f(a_i)$ is a nonzero integer with only finitely many divisors. The algorithm tries to discover g by looping over all possible sequences of integers (b_0, b_1, \dots, b_m) with b_i a divisor of $f(a_i)$ in \mathbb{Z} for all i ; note that there are only finitely many such sequences. As shown above, each sequence of this type yields exactly one $g \in \mathbb{Q}[x]$ such that $g(a_i) = b_i$ for $0 \leq i \leq m$, and we can explicitly calculate g from the a_i 's and b_i 's. Construct this finite list of g 's and see if any of them divides f , is in $\mathbb{Z}[x]$, and has $0 < \deg(g) \leq m$. If so, we have found a proper factor. If no g on this list works, then we know that no other g can possibly work, because of the equations $f(a_i) = g(a_i)h(a_i)$ that must hold in \mathbb{Z} if $f = gh$. In this latter case, the algorithm has proved irreducibility of f in $\mathbb{Q}[x]$.

For example, let us indicate how Kronecker's algorithm can be used to show mechanically that $f = x^4 + x^3 + x^2 + x + 1$ is irreducible in $\mathbb{Q}[x]$. Here $n = \deg(f) = 5$ and $m = 2$; for convenience of calculation, we choose $a_0 = 0$, $a_1 = 1$, and $a_2 = -1$. Now $f(a_0) = 1$ forces $b_0 \in \{1, -1\}$; $f(a_1) = 5$ forces $b_1 \in \{1, -1, 5, -5\}$; and $f(a_2) = 1$ forces $b_2 \in \{1, -1\}$. So there are $2 \cdot 4 \cdot 2 = 16$ sequences (b_0, b_1, b_2) to be tested. Now, if g is the unique polynomial of degree at most 2 such that $g(a_i) = b_i$ for $i = 0, 1, 2$, then $-g$ is evidently the unique polynomial of degree at most 2 such that $g(a_i) = -b_i$ for $i = 0, 1, 2$. Because of this sign symmetry, we really only have 8 sequences to test. We examine two of these sequences explicitly and let the reader (or the reader's computer) take care of the others. First, if $b_0 = b_1 = b_2 = 1$, the associated polynomial g is the constant polynomial 1. This g is, of course, a divisor of f , but we need to examine the other sequences to see if one of them yields a non-constant divisor. Second, say $b_0 = 1$, $b_1 = 5$, and $b_2 = -1$. From this data, Lagrange's interpolation formula produces

$$g = 1 \cdot \frac{(x-1)(x+1)}{(0-1)(0+1)} + 5 \cdot \frac{(x-0)(x+1)}{(1-0)(1+1)} + (-1) \cdot \frac{(x-0)(x-1)}{(-1-0)(-1-1)} = x^2 + 3x + 1.$$

Dividing f by g in $\mathbb{Q}[x]$ gives quotient $x^2 - 2x + 6$ and remainder $-15x - 5$, so that $g \nmid f$. Similarly, the g 's arising from the remaining sequences of b_i 's do not divide f either. This proves that f is, indeed, irreducible in $\mathbb{Q}[x]$.

It is nice that we can use Kronecker's algorithm to factor any specific $f \in \mathbb{Q}[x]$ or prove that f is irreducible. But, the reader will notice that Kronecker's algorithm is highly time-consuming if $\deg(f)$ is large or if the integers $f(a_i)$ have many prime divisors. Another drawback is that the algorithm cannot be used to prove irreducibility of families of polynomials that depend on one or more parameters (like $x^n - 2$ as n varies). This is one reason why it is also helpful to have general theorems such as Eisenstein's criterion that can be used to manufacture irreducible polynomials in $\mathbb{Q}[x]$ with special structure.

3.19 Algebraic Elements and Minimal Polynomials

Let F be a field, and let V be an F -algebra. Recall (§1.3) that this means V is a ring and an F -vector space such that $c(x \bullet y) = (cx) \bullet y = x \bullet (cy)$ for all $x, y \in V$ and all $c \in F$. We will fix an element $z \in V$ and study the F -algebra generated by z , which can be defined as the intersection of all subalgebras of V containing z . We write $F[\{z\}]$ to denote this subalgebra (many authors write $F[z]$ instead, but this could be confused with a polynomial ring where the indeterminate is called z). One can check as an exercise that $F[\{z\}] = \{g(z) : g \in F[x]\}$, which is the set of all elements in the algebra V that can be

written as polynomial expressions in z using coefficients from F . It follows from this that the set $S = \{1_V, z, z^2, \dots, z^n, \dots\}$ spans $F[\{z\}]$ viewed as an F -vector space.

We will obtain structural information about $W = F[\{z\}]$ in the case where W is finite-dimensional over F . In this case, the set $S = \{z^k : k \geq 0\}$ must be F -linearly dependent. So there are finitely many scalars $c_0, c_1, \dots, c_n \in F$ with $c_n \neq 0$ and $\sum_{i=0}^n c_i z^i = 0_V$. This means that z is a *root* (in V) of the nonzero polynomial $f = \sum_{i=0}^n c_i x^i \in F[x]$. By definition, we say that $z \in V$ is *algebraic over F* iff there exists a nonzero $f \in F[x]$ with $f(z) = 0$.

Now assume F is a field, V is any F -algebra, and $z \in V$ is algebraic over F . We will prove that *in the set I of all nonzero polynomials in $F[x]$ having z as a root, there exists a unique monic polynomial m of minimum degree; moreover, I consists of all nonzero multiples of m in $F[x]$.* The polynomial m will be called the *minimal polynomial of z over F* . To begin the proof, note that I is nonempty because z is algebraic over F . Since the degrees of polynomials in I constitute a nonempty subset of \mathbb{N} , we can choose a polynomial $m \in I$ of minimum degree; dividing by the leading coefficient if needed, we can arrange that m be monic. Now, let $g \in I$ be any nonzero polynomial having z as a root. Dividing g by m in $F[x]$ gives $g = mq + r$ for some $q, r \in F[x]$ with $r = 0$ or $\deg(r) < \deg(m)$. Evaluating at z gives $0 = g(z) = m(z)q(z) + r(z) = 0q(z) + r(z) = r(z)$. If r were nonzero, this would contradict the minimality of $\deg(m)$. So, in fact, $r = 0$ and $m|g$ in $F[x]$. Conversely, any multiple of m in $F[x]$ has z as a root. We have now proved everything except uniqueness of m . But if m_1 is another monic polynomial of minimum degree in I , we have just seen that $m|m_1$. Since the degrees are the same, m_1 must be a constant multiple of m . The constant must be 1, since m and m_1 are both monic. This proves that m is unique.

The minimal polynomial m of an algebraic $z \in V$ need not be irreducible in $F[x]$, but it will be irreducible if the F -algebra V is an integral domain (or, in particular, a field). For, suppose $m = fg$ for some $f, g \in F[x]$ that both have lower degree than m . Evaluating at z gives $0 = m(z) = f(z)g(z)$. If V is an integral domain, we deduce $f(z) = 0$ or $g(z) = 0$, and either possibility contradicts the minimality of the degree of m .

Once again, let z be an arbitrary element of any F -algebra V . We have seen that if $W = F[\{z\}]$ is finite-dimensional, then z is algebraic over F and has a minimal polynomial. (Note that the hypothesis will automatically hold if the whole algebra V is finite-dimensional.) Conversely, we now show that *if z is algebraic over F with minimal polynomial m of degree d , then $F[\{z\}]$ is d -dimensional with ordered basis $B = (1_V, z, z^2, \dots, z^{d-1})$.* To prove F -linear independence of B , assume $a_0, \dots, a_{d-1} \in F$ satisfy $a_0 1 + a_1 z + \dots + a_{d-1} z^{d-1} = 0_V$. Then z is a root of the polynomial $\sum_{i=0}^{d-1} a_i x^i \in F[x]$. If this polynomial were nonzero, the minimality of $\deg(m)$ would be contradicted. So this polynomial is zero, which means all a_i 's are zero. Next, to show B spans W , it suffices to show that every element in the spanning set $\{z^k : k \geq 0\}$ for W is an F -linear combination of the vectors in B . This is evidently the case for $0 \leq k < d$, since these powers of z already lie in B . Now fix $k \geq d$, and assume by induction that z^0, z^1, \dots, z^{k-1} are already known to lie in the subspace spanned by B . Write $m = x^d + \sum_{i=0}^{d-1} c_i x^i$ for some $c_i \in F$. Evaluating at z , multiplying by z^{k-d} , and solving for z^k gives $z^k = \sum_{i=0}^{d-1} -c_i z^{i+k-d}$, so that z^k is an F -linear combination of lower powers of z . All these powers already lie in the subspace spanned by B , so z^k does as well, completing the induction.

How can we find the minimal polynomial m of a specific algebraic z in a specific F -algebra V ? The key observation is that if $\deg(m) = d$, then z^d is the lowest power of z that is an F -linear combination of preceding powers of z (as we saw in the proof above). If m is not known in advance, we look at $d = 1, 2, 3, \dots$ in turn and check if the list $(1, z, z^2, \dots, z^d)$ in V is F -linearly dependent. If it is, the linear dependence relation $c_0 1 + c_1 z + c_2 z^2 + \dots + c_d z^d$ with $c_d = 1_F$ gives us the minimal polynomial $m = \sum_{i=0}^d c_i x^i$ of z .

For example, consider $z = \sqrt{2} + \sqrt{3}$, which is an element of the \mathbb{Q} -algebra \mathbb{R} . We compute

$z^0 = 1$, $z^1 = \sqrt{2} + \sqrt{3}$, $z^2 = 5 + 2\sqrt{6}$, and $z^3 = 11\sqrt{2} + 9\sqrt{3}$. These powers of z are linearly independent over \mathbb{Q} , as the reader may check. But $z^4 = 49 + 20\sqrt{6} = 10z^2 - 1$, so $x^4 - 10x^2 + 1$ must be the minimal polynomial of z over \mathbb{Q} . As the \mathbb{Q} -algebra \mathbb{R} is a field, this gives one way of proving that $x^4 - 10x^2 + 1$ must be irreducible in $\mathbb{Q}[x]$.

Now consider the F -algebra $V = M_s(F)$ of all $s \times s$ matrices with entries in F . This algebra has finite dimension (namely s^2) over F , so every matrix $A \in V$ is algebraic over F and has a minimal polynomial $m_A \in F[x]$. The degree of m_A is the dimension of the subspace of matrices spanned by $\{I_s, A, A^2, \dots, A^n, \dots\}$, where I_s denotes the $s \times s$ identity matrix. This subspace has dimension at most s^2 , so $\deg(m_A) \leq s^2$ for every $s \times s$ matrix A . For a specific example, let $s = 2$, $F = \mathbb{R}$, and $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$. The matrices I_2 and A are linearly independent, but $A^2 = \begin{bmatrix} 7 & 10 \\ 15 & 22 \end{bmatrix} = 5A + 2I_2$, so the minimal polynomial of A over \mathbb{R} is $x^2 - 5x - 2$. Observe that this polynomial is reducible in $\mathbb{R}[x]$, since it has real roots $(5 \pm \sqrt{33})/2$.

In the case of matrices A in $V = M_s(F)$, we can say a lot more about the minimal polynomial of A . The famed *Cayley–Hamilton theorem* states that A is always a root of A 's characteristic polynomial $\chi_A = \det(xI_s - A) \in F[x]$, which has degree s . Assuming this theorem, our discussion above shows that m_A must divide χ_A in $F[x]$, giving the improved bound $\deg(m_A) \leq s$. For the 2×2 matrix A considered above, m_A is equal to the characteristic polynomial $\chi_A = \det \begin{bmatrix} x-1 & -2 \\ -3 & x-4 \end{bmatrix}$, but this equality does not hold for all matrices. In later chapters, we will prove the Cayley–Hamilton theorem and use canonical forms to obtain a deeper understanding of the precise relationship between the characteristic polynomial and the minimal polynomial of a square matrix (see §5.15, §8.14, and §18.19).

3.20 Multivariable Polynomials

We conclude this chapter with a quick overview of formal power series and polynomials that involve more than one variable. Intuitively, a polynomial in variables x_1, \dots, x_m with coefficients in a given ring R is a finite sum of *monomials* of the form $cx_1^{e_1} \cdots x_m^{e_m}$, where $c \in R$ and e_1, \dots, e_m are nonnegative integers. To specify our polynomial, we must indicate which coefficient c is associated with each possible m -tuple of exponents $(e_1, \dots, e_m) \in \mathbb{N}^m$. This suggests the following precise definitions.

Given a ring R and a positive integer m , let $R[[x_1, \dots, x_m]]$ be the set of all functions from \mathbb{N}^m to R . Such a function is called a *formal power series in x_1, \dots, x_m with coefficients in R* . Given $f \in R[[x_1, \dots, x_m]]$, we introduce the formal summation notation

$$f = \sum_{e_1 \geq 0} \cdots \sum_{e_m \geq 0} f(e_1, \dots, e_m) x_1^{e_1} \cdots x_m^{e_m}$$

to represent f . We also define $R[x_1, \dots, x_m]$ to be the subset of all $f \in R[[x_1, \dots, x_m]]$ such that $f(e_1, \dots, e_m) = 0_R$ for all but finitely many inputs $(e_1, \dots, e_m) \in \mathbb{N}^m$. Elements of $R[x_1, \dots, x_m]$ are called *formal polynomials in the variables x_1, \dots, x_m* .

As in the case $m = 1$, we can define addition and multiplication operations on $R[[x_1, \dots, x_m]]$ that turn this set into a ring with subring $R[x_1, \dots, x_m]$; these rings are commutative if R is commutative. For all $f, g \in R[[x_1, \dots, x_m]]$ and all $v \in \mathbb{N}^m$, define

$(f + g)(v) = f(v) + g(v)$ and

$$(f \cdot g)(v) = \sum_{w,y \in \mathbb{N}^m: w+y=v} f(w)g(y).$$

We invite the reader to execute the tedious computations needed to verify the ring axioms. As a hint for the proof of associativity of multiplication, one should show that $(fg)h$ and $f(gh)$ are both functions from \mathbb{N}^m to R that send each $v \in \mathbb{N}^n$ to

$$\sum_{w,y,z \in \mathbb{N}^m: w+y+z=v} f(w)g(y)h(z).$$

We can identify a ring element $c \in R$ with the “constant” polynomial (or formal series) that sends $(0, 0, \dots, 0) \in \mathbb{N}^m$ to c and every other element of \mathbb{N}^m to zero. This identification allows us to view R as a subring of $R[x_1, \dots, x_m]$ or of $R[[x_1, \dots, x_m]]$. Next, define the “variable” x_i to be the function from \mathbb{N}^m to R that sends $(0, \dots, 0, 1, 0, \dots, 0)$ to 1_R (where the 1 appears in position i) and sends every other element of \mathbb{N}^m to 0_R . Using the definition of multiplication, one can then check that the ring element $cx_1^{e_1} \cdots x_m^{e_m}$ (formed by multiplying together c , then e_1 copies of x_1 , etc.) is the function from \mathbb{N}^m to R that sends (e_1, \dots, e_m) to c and everything else to zero. It follows that our formal notation

$$f = \sum_{(e_1, \dots, e_m) \in \mathbb{N}^m} f(e_1, \dots, e_m) x_1^{e_1} \cdots x_m^{e_m}$$

for the polynomial f can also be viewed as a valid algebraic identity in the ring $R[x_1, \dots, x_m]$, in which f is built up from various ring elements c and the x_i ’s by addition and multiplication in the ring. (The generalization of this comment to formal power series is not immediate, however, since we have not defined what an *infinite* sum of ring elements should mean in $F[[x_1, \dots, x_m]]$.)

As in the case $m = 1$, the multivariable polynomial ring $R[x_1, \dots, x_m]$ possesses the following universal mapping property. Suppose R and S are commutative rings, $h : R \rightarrow S$ is a ring homomorphism, and c_1, \dots, c_m are arbitrary elements of S . Then there exists a unique ring homomorphism $H : R[x_1, \dots, x_m] \rightarrow S$ such that H extends h and $H(x_i) = c_i$ for $1 \leq i \leq m$. We sketch the proof. Using the discussion in the previous paragraph, if the ring homomorphism H exists at all, then H must send $f = \sum f(e_1, \dots, e_m) x_1^{e_1} \cdots x_m^{e_m}$ to $\sum h(f(e_1, \dots, e_m)) c_1^{e_1} \cdots c_m^{e_m}$. This proves the uniqueness of H . To prove existence, take the preceding formula as the definition of H , and then check that H is a well-defined ring homomorphism extending h and sending x_i to c_i for all i .

We now state without proof some of the basic facts about divisibility for multivariable polynomials with coefficients in a field F . Given $f, g \in F[x_1, \dots, x_m]$, we say f divides g (written $f|g$, as before) iff $f = qg$ for some $q \in F[x_1, \dots, x_m]$; and we say f and g are associates iff $f = cg$ for some nonzero $c \in F$. A polynomial $p \in F[x_1, \dots, x_m]$ is irreducible in $F[x_1, \dots, x_m]$ iff the only divisors of p in $F[x_1, \dots, x_m]$ are nonzero constants and associates of p . One can prove that every nonzero $f \in F[x_1, \dots, x_m]$ can be factored uniquely into the form $c_f \prod_{i \in I} p_i^{e_i(f)}$, where $c_f \in F$ is nonzero, the p_i ’s range over the set of all monic irreducible polynomials in $F[x_1, \dots, x_m]$, and each $e_i(f) \in \mathbb{N}$. Then divisibility, gcds, and lcms of multivariable polynomials can be analyzed in terms of these irreducible factorizations, as we did in §3.14 for $m = 1$. However, we warn the reader that some techniques and results used in the one-variable case do not extend to polynomials in more than one variable. For instance, given $f, g \in F[x_1, \dots, x_m]$ with $m > 1$, it need not be true that $d = \gcd(f, g)$ is a polynomial combination of f and g . For a specific example, one may check that $\gcd(x_1, x_2) = 1$, but 1 cannot be written as $ax_1 + bx_2$ for any choice of

$a, b \in F[x_1, \dots, x_m]$. There is a division algorithm for multivariable polynomials, but it is much more subtle than the one-variable version. We will not go into the details here, but instead refer the reader to the outstanding book by Cox, Little, and O'Shea [11].

3.21 Summary

1. *Polynomials and Formal Power Series.* Informally, a polynomial with coefficients in a ring R is an expression $\sum_{n=0}^d f_n x^n$ with $d \in \mathbb{N}$ and each $f_n \in R$, whereas a formal power series is an expression $\sum_{n=0}^{\infty} f_n x^n$ with each $f_n \in R$. More precisely, the set $R[[x]]$ of formal power series consists of all sequences $f = (f_n : n \in \mathbb{N})$ with each $f_n \in R$. $R[[x]]$ is a ring under the operations

$$(f_n : n \in \mathbb{N}) + (g_n : n \in \mathbb{N}) = (f_n + g_n : n \in \mathbb{N}),$$

$$(f_n : n \in \mathbb{N}) \cdot (g_n : n \in \mathbb{N}) = \left(\sum_{j, k \in \mathbb{N}: j+k=n} f_j g_k : n \in \mathbb{N} \right).$$

The set $R[x]$ of formal polynomials consists of all $(f_n : n \in \mathbb{N})$ in $R[[x]]$ such that $f_n = 0$ for all but finitely many $n \in \mathbb{N}$. $R[x]$ is a subring of $R[[x]]$ containing an isomorphic copy of R (the constant polynomials). R is commutative iff $R[x]$ is commutative, and R is an integral domain iff $R[x]$ is an integral domain; similarly for $R[[x]]$. $R[x]$ is never a field. If F is a field, then $F[x]$ and $F[[x]]$ are infinite-dimensional F -vector spaces and F -algebras. There are similar definitions and results for formal power series and polynomials in m variables x_1, \dots, x_m , obtained by replacing \mathbb{N} by \mathbb{N}^m in the formulas above.

2. *Degree.* The degree of a nonzero polynomial $f = \sum_{n \geq 0} f_n x^n$ is the largest $n \in \mathbb{N}$ with $f_n \neq 0$. The degree of the zero polynomial is undefined. For nonzero polynomials p, q such that $p + q \neq 0$, $\deg(p + q) \leq \max(\deg(p), \deg(q))$, and strict inequality occurs iff $\deg(p) = \deg(q)$ and the leading terms of p and q sum to zero. Similarly, if the leading terms of p and q do not multiply to zero, then the degree addition formula $\deg(pq) = \deg(p) + \deg(q)$ is valid.
3. *Polynomial Functions, Evaluation Homomorphisms, and the UMP for $R[x]$.* Given a ring R and $p = \sum_{i=0}^n a_i x^i \in R[x]$, the associated polynomial function $f_p : R \rightarrow R$ is given by $f_p(c) = p(c) = \sum_{i=0}^n a_i c^i$ for $c \in R$. Polynomial rings satisfy the following universal mapping property (UMP): for any commutative rings R and S and any ring homomorphism $h : R \rightarrow S$ and any $c \in S$, there exists a unique ring homomorphism $H : R[x] \rightarrow S$ extending h such that $H(x) = c$. The special case $h = \text{id}_R : R \rightarrow R$ yields the evaluation homomorphisms $E_c : R[x] \rightarrow R$ given by $E_c(p) = f_p(c) = p(c)$ for all $p \in R[x]$.
4. *Divisibility Definitions.* For a field F and $f, g \in F[x]$, f divides g in $F[x]$ (written $f|g$) iff $g = qf$ for some $q \in F[x]$. Polynomials $f, g \in F[x]$ are associates iff $g = cf$ for some nonzero $c \in F$; every nonzero polynomial has a unique monic associate. We say $f \in F[x]$ is a greatest common divisor (gcd) of nonzero $g_1, \dots, g_k \in F[x]$ iff f divides all g_i and has maximum degree among all common divisors of the g_i 's; $\gcd(g_1, \dots, g_k)$ refers to the unique monic gcd of the g_i 's. Similarly, $h \in F[x]$ is a least common multiple (lcm) of the g_i 's iff every g_i divides h and h has minimum degree among all common multiples of the g_i 's; $\text{lcm}(g_1, \dots, g_k)$ is the

unique monic lcm of the g_i 's. A polynomial $p \in F[x]$ is irreducible in $F[x]$ iff its only divisors are constants and associates of p . Whether p is irreducible depends on the field F of coefficients.

5. *Polynomial Division with Remainder.* For a field F and $f, g \in F[x]$ with $g \neq 0$, there exist a unique quotient $q \in F[x]$ and remainder $r \in F[x]$ with $f = qg + r$ and either $r = 0$ or $\deg(r) < \deg(g)$. There is an algorithm to compute q and r from f and g . Replacing F by a commutative ring R , q and r will still exist if the leading coefficient of g is invertible in R ; q and r will still be unique if R is an integral domain.
6. *Greatest Common Divisors.* Given a field F and nonzero $f_1, \dots, f_n \in F[x]$, there exists a unique monic gcd g of the f_i 's, and the common divisors of the f_i 's are precisely the divisors of g . There exist $u_1, \dots, u_n \in F[x]$ with $g = u_1f_1 + \dots + u_nf_n$; this fact does not extend to polynomials in more than one variable. When $n = 2$, we can find g by repeatedly dividing the previous remainder by the current remainder (starting with f_1 divided by f_2), letting g be the last nonzero remainder, then working backwards to find u_1 and u_2 . For any $n \geq 1$, we can find g and the u_i by starting with an $n \times n$ identity matrix augmented with a column containing f_1, \dots, f_n , row-reducing this matrix until a single nonzero entry g remains in the extra column, and looking at the other entries in g 's row to find u_1, \dots, u_n .
7. *Roots of Polynomials.* Given a field F , $p \in F[x]$, and $c \in F$, c is a root of p (i.e., $p(c) = 0$) iff $(x - c)|p$ in $F[x]$. If $\deg(p) = n$, then p has at most n distinct roots in F . For finite F , we can find the roots of p in F by exhaustive search. For $F = \mathbb{Q}$, we can replace p by an associate $\sum_{i=0}^n a_i x^i$ in $\mathbb{Z}[x]$; all rational roots r/s of p must satisfy $r|a_0$ and $s|a_n$.
8. *Irreducible Polynomials.* For a field F , every $p \in F[x]$ of degree 1 is irreducible. Given $p \in F[x]$ with $\deg(p) > 1$, p irreducible implies p has no root in F ; the converse is valid only for $\deg(p) = 2$ or 3. For F finite, one can factor $p \in F[x]$ or prove p is irreducible by dividing by finitely many potential divisors of degree at most $\deg(p)/2$. For $F = \mathbb{C}$, p is irreducible iff $\deg(p) = 1$. For $F = \mathbb{R}$, p is irreducible iff $\deg(p) = 1$ or $p = ax^2 + bx + c$ with $a, b, c \in \mathbb{R}$, $a \neq 0$, and $b^2 - 4ac < 0$. If an irreducible $p \in F[x]$ divides a product $f_1 f_2 \cdots f_n$ with all $f_i \in F[x]$, then $p|f_i$ for some i ; this property characterizes irreducible polynomials.
9. *Unique Factorization in $F[x]$.* For F a field, every non-constant $g \in F[x]$ can be written $g = cp_1 p_2 \cdots p_r$ where $0 \neq c \in F$ and each p_i is monic and irreducible in $F[x]$; this factorization is unique except for reordering the p_i 's. We can also write the factorization as $g = c(g) \prod_{i \in I} p_i^{e_i(g)}$ where $c(g) \in F$ is nonzero, $\{p_i : i \in I\}$ is an indexed set of all monic irreducible polynomials in $F[x]$, and each $e_i(g) \in \mathbb{N}$. For $f, g \in F[x]$, $f|g$ iff $e_i(f) \leq e_i(g)$ for all $i \in I$; it follows that we can compute gcds (resp. lcms) by the formulas $e_i(\gcd(f_1, \dots, f_n)) = \min(e_i(f_1), \dots, e_i(f_n))$ and $e_i(\text{lcm}(f_1, \dots, f_n)) = \max(e_i(f_1), \dots, e_i(f_n))$ for all $i \in I$. In particular, $\text{lcm}(f, g) \gcd(f, g) = fg$.
10. *Criteria for Irreducibility in $\mathbb{Q}[x]$.* If $u \in \mathbb{Z}[x]$ has degree $n > 0$ and is reducible in $\mathbb{Q}[x]$, then $u = gh$ for some $g, h \in \mathbb{Z}[x]$ with $\deg(g), \deg(h) < n$. Given $f = \sum_{i=0}^n f_i x^i \in \mathbb{Z}[x]$ of degree n , suppose the reduction of f mod p is irreducible in $\mathbb{Z}_p[x]$ for some prime integer p not dividing f_n ; then f is irreducible in $\mathbb{Q}[x]$. If there is a prime $p \in \mathbb{Z}$ such that $p|f_i$ for all $i < n$, p does not divide f_n , and $p^2 \nmid f_0$, then f is irreducible in $\mathbb{Q}[x]$ (Eisenstein's criterion).
11. *Lagrange Interpolation Formula.* Given a field F , distinct $a_0, \dots, a_n \in F$, and

arbitrary $b_0, \dots, b_n \in F$, there exists a unique $g \in F[x]$ with $g = 0$ or $\deg(g) \leq n$ such that $g(a_i) = b_i$ for all i . Explicitly,

$$g = \sum_{i=0}^n b_i p_i, \text{ where } p_i = \prod_{j \neq i} (x - a_j) \quad \left/ \prod_{j \neq i} (a_i - a_j) \text{ for all } i. \right.$$

12. *Kronecker's Factoring Algorithm in $\mathbb{Q}[x]$.* Given $f \in \mathbb{Z}[x]$ of degree n , Kronecker's algorithm produces a finite list of potential divisors of f of degree at most $m = \lfloor n/2 \rfloor$ as follows. Pick distinct $a_0, \dots, a_m \in \mathbb{Z}$; loop over all $(b_0, \dots, b_m) \in \mathbb{Z}^{m+1}$ such that $b_i | f(a_i)$ for all i ; use Lagrange interpolation to build $g \in \mathbb{Q}[x]$ with $g(a_i) = b_i$ for all i . If one of these g 's is non-constant in $\mathbb{Z}[x]$ and divides f , we have partially factored f ; if no g works, then f must be irreducible in $\mathbb{Q}[x]$. As a special case, if some $f(a_i) = 0$, we know $(x - a_i) | f$ in $\mathbb{Q}[x]$.
13. *Minimal Polynomials.* Let F be a field, V an F -algebra, and $z \in V$. We say z is algebraic over F iff $g(z) = 0$ for some nonzero $g \in F[x]$. For algebraic z , there exists a unique monic polynomial $m \in F[x]$ with $m(z) = 0$ such that m divides every other polynomial in $F[x]$ having z as a root; m is called the minimal polynomial of z over F . The subalgebra $F[\{z\}]$ spanned by all powers of z is finite-dimensional iff z is algebraic over F ; in this case, a basis for this subalgebra (viewed as an F -vector space) is $(1, z, z^2, \dots, z^{d-1})$ where d is the degree of z 's minimal polynomial. The minimal polynomial of z must be irreducible for an integral domain or field V , but can be reducible in other algebras. By the Cayley–Hamilton theorem, every $s \times s$ matrix $A \in M_s(F)$ has a minimal polynomial of degree at most s that divides the characteristic polynomial of A .

3.22 Exercises

1. Let R be a ring. (a) Prove carefully that $R[[x]]$ is a commutative group under addition, indicating which ring axioms for R are needed at each point in the proof. (b) Prove that $R[x]$ is an additive subgroup of $R[[x]]$. (c) Repeat (a) and (b) for $R[[x_1, \dots, x_m]]$ and $R[x_1, \dots, x_m]$.
2. Let R be a ring. (a) Carefully prove the ring axioms for $R[[x]]$ involving multiplication, indicating which ring axioms for R are needed at each point in the proof. Prove $R[[x]]$ is commutative iff R is commutative. (b) Repeat (a) for $R[[x_1, \dots, x_m]]$. (c) Prove $R[x_1, \dots, x_m]$ is a subring of $R[[x_1, \dots, x_m]]$.
3. Let F be a field. (a) Prove that $F[[x]]$ is an F -vector space and F -algebra using the scalar multiplication $c \cdot (f_i : i \geq 0) = (cf_i : i \geq 0)$ for $c, f_i \in F$. (b) Prove that $F[x]$ is a subspace of $F[[x]]$. (c) Prove that $\{1, x, x^2, \dots, x^n, \dots\}$ is an F -linearly independent subset of $F[x]$. (d) Repeat (a) and (b) for formal series and polynomials in m variables, and show that $\{x_1^{e_1} x_2^{e_2} \cdots x_m^{e_m} : (e_1, \dots, e_m) \in \mathbb{N}^m\}$ is a basis of $F[x_1, \dots, x_m]$.
4. Let R be a ring. (a) Show that $j : R \rightarrow R[x]$ given by $j(a) = (a, 0, 0, \dots)$ for $a \in R$ is an injective ring homomorphism. (b) For $m \geq 2$, define a similar map $j_m : R \rightarrow R[x_1, \dots, x_m]$ and verify that it is an injective ring homomorphism.
5. Let R be a ring with subring S and ideal I . (a) Prove that $S[[x]]$ is a subring of $R[[x]]$, and $S[x]$ is a subring of $R[x]$. (b) Write $I[[x]]$ (resp. $I[x]$) for the subset of

formal series (resp. polynomials) all of whose coefficients lie in I . Prove there are ring isomorphisms $R[[x]]/I[[x]] \cong (R/I)[[x]]$ and $R[x]/I[x] \cong (R/I)[x]$.

6. Let R be a ring. Let $x = (0_R, 1_R, 0_R, 0_R, \dots) \in R[x]$. (a) Prove carefully that for $i \geq 1$, the product of i copies of x in $R[x]$ is the sequence $(e_j : j \geq 0)$ with $e_i = 1$ and $e_j = 0$ for all $j \neq i$. (b) For $m > 1$ and $1 \leq k \leq m$, define $x_k : \mathbb{N}^m \rightarrow R$ by letting x_k map $(0, \dots, 1, \dots, 0)$ to 1_R (the 1 is in position k) and letting x_k map everything else in \mathbb{N}^m to 0_R . Prove carefully that $x_1^{e_1} x_2^{e_2} \cdots x_m^{e_m}$ (the product of $e_1 + \cdots + e_m$ elements of $R[x_1, \dots, x_m]$) is the function mapping (e_1, \dots, e_m) to 1_R and everything else in \mathbb{N}^m to 0_R . (This justifies the monomial notation used for multivariable formal series and polynomials.)

7. *Binomial Theorem.* Let R be a ring. (a) Prove: for all $r \in R$ and $n \in \mathbb{N}^+$, the identity

$$(r+x)^n = \sum_{k=0}^n \binom{n}{k} r^{n-k} x^k$$

holds in $R[x]$. Here, $\binom{n}{k} = n!/(k!(n-k)!)$, and the notation js (for $j \in \mathbb{N}$, $s \in R$) denotes the sum of j copies of s in R . (b) Give a specific example showing that the identity in (a) can be false if the polynomial x is replaced by an element of R .

8. Prove: for all $m > 1$ and all rings R , there are ring isomorphisms $R[[x_1, \dots, x_{m-1}]][[x_m]] \cong R[[x_1, \dots, x_{m-1}, x_m]]$ and $R[x_1, \dots, x_{m-1}][x_m] \cong R[x_1, \dots, x_m]$. (This gives an alternative recursive definition of multivariable formal series and polynomials.)
9. Let R be a ring. Define the *order* of a nonzero formal series $f \in R[[x]]$, denoted $\text{ord}(f)$, to be the smallest $n \in \mathbb{N}$ with $f_n \neq 0_R$. (a) Suppose f, g are nonzero formal series with $f + g \neq 0$. Find and prove an inequality relating $\text{ord}(f+g)$ to $\text{ord}(f)$ and $\text{ord}(g)$, and determine when strict inequality holds. (b) Under what conditions on $f, g \in R[[x]]$ will it be true that $\text{ord}(fg) = \text{ord}(f) + \text{ord}(g)$? Prove your answer.
10. (a) Prove that a ring R is an integral domain iff $R[[x]]$ is an integral domain. (b) Fix $m > 1$. Prove R is an integral domain iff $R[x_1, \dots, x_m]$ is an integral domain iff $R[[x_1, \dots, x_m]]$ is an integral domain.
11. Give a specific example of a commutative ring R and polynomials of all degrees in $R[x]$ that have multiplicative inverses in $R[x]$.
12. Let F be a field. Prove that $g \in F[[x]]$ is invertible in $F[[x]]$ iff the constant coefficient in g is nonzero. Is $F[[x]]$ ever a field?
13. Let F be a field and $n \in \mathbb{N}^+$. Is the set of $p \in F[x]$ with $\deg(p) = n$ or $p = 0$ a subspace of $F[x]$? What about the set of p with $\deg(p) > n$ or $p = 0$? What about the set of $p \in F[[x]]$ with $\text{ord}(p) > n$ or $p = 0$? (See Exercise 9 for the definition of $\text{ord}(p)$.)
14. Suppose F is a field and, for all $n \in \mathbb{N}$, $p_n \in F[x]$ satisfies $\deg(p_n) = n$. Prove $\{p_n : n \in \mathbb{N}\}$ is a basis of the F -vector space $F[x]$.
15. Let $p = x^4 + 3x^2 + 4x + 1 \in \mathbb{Z}_5[x]$. Find the associated polynomial function $f_p : \mathbb{Z}_5 \rightarrow \mathbb{Z}_5$.
16. (a) Explicitly describe all functions $g : \mathbb{Z}_2 \rightarrow \mathbb{Z}_2$. (b) For each g found in (a), find *all* polynomials $p \in \mathbb{Z}_2[x]$ such that $f_p = g$.

17. Let F be a finite field. (a) Prove that for every function $g : F \rightarrow F$, there exists at least one $p \in F[x]$ such that $g = f_p$. (b) Prove or disprove: for each $g : F \rightarrow F$, there exist infinitely many $p \in F[x]$ with $g = f_p$.
18. Give a specific example of an infinite field F and a function $g : F \rightarrow F$ such that $g \neq f_p$ for all polynomials $p \in F[x]$.
19. Given a field F , let S be the set of all functions from F to F . (a) Verify that S is a commutative ring under pointwise operations on functions $((f+g)(c) = f(c)+g(c)$ and $(f \cdot g)(c) = f(c) \cdot g(c)$ for all $f, g \in S$ and all $c \in F$). (b) Define $\phi : F[x] \rightarrow S$ by $\phi(p) = f_p$, where f_p is the polynomial function sending $c \in F$ to $p(c)$. Prove carefully that ϕ is a ring homomorphism. (c) Prove: ϕ is surjective iff F is finite (for infinite F , you may need a cardinality argument). (d) Prove: ϕ is injective iff F is infinite. (Part (d) shows that it is only safe to “identify” formal polynomials with polynomial functions when the field of coefficients is infinite.)
20. Let R be a commutative ring and $c \in R$. Prove that the map $\phi_c : R[x] \rightarrow R[x]$ given by $\phi_c(p) = p(x - c)$ for $p \in R[x]$ is a ring isomorphism. (Use the UMP.)
21. Give two proofs of the universal mapping property for $R[x_1, \dots, x_m]$ stated in §3.20, by: (a) filling in the details of the proof sketched in the text; (b) using induction on m , the UMP for one-variable polynomial rings, and the result of Exercise 8.
22. Suppose R and S are commutative rings and $f : R \rightarrow S$ is a ring homomorphism. (a) Use the UMP to show that $F : R[x] \rightarrow S[x]$ defined by $F(\sum_{i \geq 0} a_i x^i) = \sum_{i \geq 0} f(a_i)x^i$ is a ring homomorphism. (b) Prove: if f is one-to-one (resp. onto), then F is one-to-one (resp. onto).
23. (a) Let R be a subring of a commutative ring S , and let $c \in S$. Prove that $R[\{c\}] = \{f_p(c) : p \in R[x]\}$ is the smallest subring of S containing R and c . (b) Let F be a field and V a (not necessarily commutative) F -algebra. For $z \in V$, prove that $F[\{z\}] = \{f_p(z) : p \in F[x]\}$ is an F -subalgebra of V . Prove also that $F[\{z\}]$ is spanned (as an F -vector space) by $\{z^n : n \geq 0\}$ and is the smallest subalgebra of V containing F and z .
24. Let $h : M_2(\mathbb{R}) \rightarrow M_2(\mathbb{R})$ be the identity map. Show there does not exist any extension of h to a ring homomorphism $H : M_2(\mathbb{R})[x] \rightarrow M_2(\mathbb{R})$ such that $H(x) = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$. (This shows that the UMP for polynomial rings can fail if the coefficient ring is not commutative.)
25. For each $f, g \in \mathbb{Q}[x]$, find the quotient q and remainder r when f is divided by g .
 (a) $f = x^8 + x^4 + 1$, $g = x^4 - x^2 + 1$; (b) $f = 3x^3 - 5x^2 + 11x$, $g = x^2 - 2x + 3$;
 (c) $f = x^4 - 5x^2 - x + 1$, $g = x - 2$; (d) $f = x^5 - 1$, $g = 2x^2 - 1$.
26. For each field F and $f, g \in F[x]$, find the quotient q and remainder r when f is divided by g . (a) $F = \mathbb{Z}_2$, $f = x^5 + x^4 + x + 1$, $g = x^3 + x$; (b) $F = \mathbb{Z}_3$, $f = x^5 + x^4 + x + 1$, $g = x^3 + x$; (c) $F = \mathbb{Z}_5$, $f = 3x^4 + x^3 + 2x + 4$, $g = 2x^2 + 3x + 1$;
 (d) $F = \mathbb{Z}_7$, $f = 3x^4 + x^3 + 2x + 4$, $g = 2x^2 + 3x + 1$.
27. (a) Give a specific example of an integral domain R and nonzero $f, g \in R[x]$ such that there do not exist $q, r \in R[x]$ satisfying $f = qg + r$ and $r = 0$ or $\deg(r) < \deg(g)$. (b) Give a specific example of a commutative ring R and $f, g \in R[x]$ such that the quotient q when f is divided by g (with remainder) is not unique.
 (c) Give a specific example of a commutative ring R and $f, g \in R[x]$ such that the remainder r when f is divided by g is not unique.

28. *Division Theorem for \mathbb{Z} .* Prove that for all $a, b \in \mathbb{Z}$ with b nonzero, there exist unique $q, r \in \mathbb{Z}$ with $a = bq + r$ and $0 \leq r < |b|$.
29. *Euclid's Algorithm for Integer GCDs.* Prove: for all nonzero integers a_1, \dots, a_n , there exists a unique positive gcd d of a_1, \dots, a_n , and there exist integers u_1, \dots, u_n with $d = u_1 a_1 + \dots + u_n a_n$. Moreover, the common divisors of the a_i 's are precisely the divisors of d . For $n = 2$, describe an algorithm (involving repeated integer division) for finding d and u_1, u_2 from a_1, a_2 .
30. Prove: for all integers p, r, s , if p is prime and $p|r s$, then $p|r$ or $p|s$.
31. *Fundamental Theorem of Arithmetic.* Prove that every nonzero integer n can be factored as $n = c p_1 p_2 \cdots p_k$, where $c \in \{1, -1\}$, $k \geq 0$, and each p_i is a prime positive integer. Also show that the factorization is unique except for reordering the p_i 's.
32. *Criterion for Invertibility mod n .* Prove: for all $n, t \in \mathbb{N}^+$, $\gcd(n, t) = 1$ iff there exists $s \in \mathbb{Z}_n$ with $st \equiv 1 \pmod{n}$; and s is unique when it exists.
33. Let F be a field and $f, g, h \in F[x]$. (a) Prove: if $f|gh$ and $\gcd(f, g) = 1$, then $f|h$. (b) Prove: if $\gcd(f, g) = 1$ and $f|h$ and $g|h$, then $fg|h$. (c) Prove: if $d \in F[x]$ is a gcd of f and g , then hd is a gcd of hf and hg . (d) Give examples to show that (a) and (b) can fail without the hypothesis $\gcd(f, g) = 1$. (e) Repeat (a) through (d) for integers f, g, h, d .
34. Let F be a field, and assume $f, f_1 \in F[x]$ are associates. (a) Prove: for all $h \in F[x]$, $f|h$ iff $f_1|h$; and $h|f$ iff $h|f_1$. (b) Prove: f is irreducible iff f_1 is irreducible.
35. Let R be any commutative ring. For $r, s \in R$, define $r|s$ iff $s = qr$ for some $q \in R$. Define $r, s \in R$ to be *associates in R* , denoted $r \sim s$, iff $r|s$ and $s|r$. (a) Prove that the divisibility relation $|$ on R is reflexive and transitive. (b) Prove that the association relation \sim is an equivalence relation on R . (c) Prove that for all integral domains R and all $r, s \in R$, $r \sim s$ iff $s = ur$ for some invertible element $u \in R$. (d) Describe the equivalence classes of \sim for these rings: $R = \mathbb{Z}$; R is a field; $R = F[x]$, where F is a field; $R = \mathbb{Z}_{12}$.
36. For each $f, g \in \mathbb{Q}[x]$, use Euclid's algorithm to compute $d = \gcd(f, g)$ and $u, v \in \mathbb{Q}[x]$ with $d = uf + vg$. (a) $f = x^6 - 1$, $g = x^4 - 1$; (b) $f = x^3 - 5$, $g = x^2 - 2$; (c) $f = x^3 + x + 1$, $g = 2x^2 - 3x + 1$; (d) $f = x^5 + x^4 - 3x^2 - 2x - 2$, $g = x^4 - x^3 - x^2 + 6$.
37. Let $f_1 = x^2 - 3x + 2$, $f_2 = x^2 - 5x + 6$, and $f_3 = x^2 - 4x + 3$ in $\mathbb{Q}[x]$. Follow the proof in §3.9 to find $d = \gcd(f_1, f_2, f_3)$ and $u_1, u_2, u_3 \in \mathbb{Q}[x]$ with $d = u_1 f_1 + u_2 f_2 + u_3 f_3$.
38. Let K be a field with subfield F . Prove that for all nonzero $f_1, \dots, f_n \in F[x]$, the gcd of the f_i 's in the ring $F[x]$ equals the gcd of the f_i 's in the ring $K[x]$ (although the sets of common divisors of the f_i 's and the irreducible factorizations of the f_i 's are different in these two rings).
39. Let F be a field. (a) Use the division theorem to prove that for every ideal I of $F[x]$, there exists $g \in F[x]$ such that $I = F[x]g = \{hg : h \in F[x]\}$. (This says that $F[x]$ is a *principal ideal domain*, or *PID*. We will study PIDs in depth in Chapter 18.) (b) Given ideals $I = F[x]g$ and $J = F[x]h$ in $F[x]$ (where $g, h \in F[x]$ are generators of the ideals), show that $I \subseteq J$ iff $h|g$ in $F[x]$. Deduce that $I = J$ iff h and g are associates in $F[x]$. (c) Given nonzero $f_1, \dots, f_n \in F[x]$, reprove the theorem that these polynomials have a unique monic gcd g , and that g has the form $\sum_{i=1}^n u_i f_i$ for some $u_i \in F[x]$, by studying the ideal $J = \{u_1 f_1 + \dots + u_n f_n : u_i \in F[x]\}$. (d) Given nonzero $f_1, \dots, f_n \in F[x]$, the set $K = \bigcap_{i=1}^n F[x]f_i$ is an ideal of $F[x]$, so is generated by some $h \in F[x]$. Find (with proof) the specific relation between h and the f_i 's.

40. For each field F and $f, g \in F[x]$, use the matrix reduction algorithm in §3.10 to find $d = \gcd(f, g)$ and $u, v \in F[x]$ with $d = uf + vg$. (a) $F = \mathbb{Q}$, $f = x^{12} - 1$, $g = x^8 - 1$; (b) $F = \mathbb{Z}_5$, $f = x^3 + x^2 + 3x + 2$, $g = 2x^4 + 2x^2 + x + 3$; (c) $F = \mathbb{Z}_3$, $f = x^3 + 2x^2 + 2$, $g = x^3 + x^2 + 2x + 2$; (d) $F = \mathbb{Z}_7$, $f = x^3 + 2x^2 + 4x + 1$, $g = x^2 + 5x + 6$; (e) $F = \mathbb{Z}_2$, $f = x^8 + x^4 + x^3 + x + 1$, $g = x^5 + x^2 + x$.
41. Give a careful proof that the matrix reduction algorithm in §3.10 for finding the gcd of a list of n polynomials terminates and gives correct results.
42. Let $f_1 = x^6 + x^4 + x + 1$, $f_2 = x^8 + x^6 + x^5 + x^4 + x^3 + x^2 + 1$, and $f_3 = x^6 + x^5 + x^2 + 1$ in $\mathbb{Z}_2[x]$. Use matrix reduction to find $d = \gcd(f_1, f_2, f_3)$ and $u_1, u_2, u_3 \in \mathbb{Z}_2[x]$ with $d = u_1f_1 + u_2f_2 + u_3f_3$.
43. *Matrix Reduction Algorithm for Integer GCDs.* Describe a modification of the matrix reduction algorithm in §3.10 that takes as input nonzero integers a_1, \dots, a_n and returns as output $d = \gcd(a_1, \dots, a_n)$ and integers u_1, \dots, u_n with $d = \sum_{i=1}^n u_i a_i$. Prove that your algorithm terminates and returns a correct answer.
44. For each a, b below, use the algorithm in the previous problem to compute $d = \gcd(a, b)$ and find integers u, v with $d = au + bv$. (a) $a = 101$, $b = 57$; (b) $a = 516$, $b = 215$; (c) $a = 1300$, $b = 967$; (d) $a = 1702$, $b = 483$.
45. Let F be a field. Reprove the theorem that a polynomial $p \in F[x]$ of degree n has at most n roots in F by appealing to the uniqueness of the irreducible factorization of p .
46. Prove: for all $n \in \mathbb{N}^+$, there exists a commutative ring R such that $x^2 - 1_R \in R[x]$ has more than n roots in R .
47. Given $f = x^5 + 6x^4 + 9x^3 + 3x^2 + 10x + 11 \in \mathbb{Z}_{13}[x]$, find all roots of f in \mathbb{Z}_{13} .
48. Prove: for any field F and any finite subset S of F , there exist infinitely many $p \in F[x]$ such that S is the set of roots of p in F .
49. Prove: for any finite field F of size n , $\prod_{c \in F} (x - c) = x^n - x$ in $F[x]$. [Hint: use Exercise 8(b) in Chapter 1.]
50. Given a polynomial $f = 7x^5 + 10x^4 + \dots - 9 \in \mathbb{Z}[x]$ of degree 5 (where the middle coefficients are unknown integers), list all possible roots of f in \mathbb{Q} .
51. Find all rational roots of $f = 10x^5 + 11x^4 - 41x^3 + 29x^2 - 51x + 18$, and hence factor f into irreducible polynomials in $\mathbb{Q}[x]$ and in $\mathbb{C}[x]$.
52. Find all $x \in \mathbb{Q}$ solving $0 = 6x^6 + 18.6x^5 + 12.6x^4 - 5.4x^3 - 25.2x^2 + 24.6x - 6$.
53. Use the fundamental theorem of algebra to prove the assertion in the text that $f \in \mathbb{R}[x]$ is irreducible iff $\deg(f) = 1$ or $f = ax^2 + bx + c$ for some $a, b, c \in \mathbb{R}$ with $a \neq 0$ and $b^2 - 4ac < 0$.
54. Find all irreducible polynomials of degree 5 or less in $\mathbb{Z}_2[x]$. Explain how you know your answers are irreducible, and how you know you have found all of them.
55. (a) Find all monic irreducible polynomials of degree two in $\mathbb{Z}_3[x]$. Explain how you know your answers are irreducible, and how you know you have found all of them. (b) For any prime p , count the number of monic irreducible degree 2 polynomials in $\mathbb{Z}_p[x]$.
56. (a) Prove that $x^4 + x^2 + 2x + 1$ is irreducible in $\mathbb{Z}_3[x]$. (b) Describe all polynomials in $\mathbb{Q}[x]$ whose irreducibility in $\mathbb{Q}[x]$ can be deduced via part (a) and reduction mod 3. (c) Can you find a way to use (a) to deduce the irreducibility of $x^4 + 8x^3 + 25x^2 + 38x + 25$ in $\mathbb{Q}[x]$?

57. For each field F , decide (with proof) whether $x^4 + x^3 + x^2 + x + 1$ is irreducible in $F[x]$: (a) $F = \mathbb{C}$; (b) $F = \mathbb{R}$; (c) $F = \mathbb{Q}$; (d) $F = \mathbb{Z}_2$; (e) $F = \mathbb{Z}_3$; (f) $F = \mathbb{Z}_5$.
58. For each field F , decide (with proof) whether $(1_F, 0_F, -1_F, 0_F, 1_F, 0_F, 0_F, \dots)$ is irreducible in $F[x]$: (a) $F = \mathbb{Z}_2$; (b) $F = \mathbb{Z}_3$; (c) $F = \mathbb{Z}_5$; (d) $F = \mathbb{Z}_7$; (e) $F = \mathbb{Z}_{13}$; (f) $F = \mathbb{Q}$; (g) $F = \mathbb{C}$.
59. Prove that $f = x^4 + 3x^2 + 1$ is irreducible in $\mathbb{Q}[x]$ as follows. (a) Prove f has no rational root. (b) Explain why the reducibility of f would guarantee the existence of monic quadratics $g, h \in \mathbb{Z}[x]$ with $f = gh$. (c) Write $g = x^2 + bx + c$, $h = x^2 + dx + e$ for $b, c, d, e \in \mathbb{Z}$. By comparing the coefficients of gh to the coefficients of f , show that g, h as in (b) do not exist.
60. Prove that each polynomial below is irreducible in $\mathbb{Q}[x]$, using the indicated methods. (a) $x^5 - 4x^2 + 2$ (use reduction mod a prime); (b) $x^7 + 30x^5 - 6x^2 + 12$ (use Eisenstein's criterion); (c) $x^4 - 5x^2 + 1$ (use any method); (d) $x^3 + ax + 1$, where $a \in \mathbb{Z}$ is a fixed integer not equal to -2 or 0 (use any method).
61. It can be shown that $x^7 + x^3 + 1$ is irreducible in $\mathbb{Z}_2[x]$. Using any method, prove the irreducibility of the following polynomials: (a) $x^5 + 30x^4 + 210x^3 + 300$ in $\mathbb{Q}[x]$; (b) $x^3 + 2x + 1$ in $\mathbb{Z}_5[x]$; (c) $x^7 + 20x^5 - 8x^4 + 11x^3 + 14x - 9$ in $\mathbb{Q}[x]$; (d) $x^4 - 8$ in $\mathbb{Q}[x]$.
62. Find the factorization of $x^8 - 1$ into a product of irreducible polynomials in each of the following polynomial rings: (a) $\mathbb{C}[x]$; (b) $\mathbb{R}[x]$; (c) $\mathbb{Q}[x]$; (d) $\mathbb{Z}_2[x]$. State briefly how you know that the factors are irreducible.
63. Find the factorization of $x^{12} - 1$ into a product of irreducible polynomials in each of the following polynomial rings: (a) $\mathbb{C}[x]$; (b) $\mathbb{R}[x]$; (c) $\mathbb{Q}[x]$; (d) $\mathbb{Z}_2[x]$; (e) $\mathbb{Z}_3[x]$; (f) $\mathbb{Z}_5[x]$. In each case, explain how you know that the factors are irreducible.
64. (a) Use the rational root theorem to prove that for all $a, n \in \mathbb{N}^+$, $\sqrt[n]{a}$ is rational iff $a = b^n$ for some integer b . (b) Reprove (a) using the uniqueness of prime factorizations of integers. (c) Deduce that $\sqrt[n]{a}$ is rational iff $a = \prod_i p_i^{e_i}$ for some distinct primes p_i and integers $e_i \geq 0$ where $n|e_i$ for all i .
65. Let F be a field, and assume $p = a_0 + a_1x + a_2x^2 + \dots + a_{n-1}x^{n-1} + a_nx^n$ is irreducible in $F[x]$, where $a_i \in F$ and $a_n \neq 0$. Let $q = a_n + a_{n-1}x + a_{n-2}x^2 + \dots + a_1x^{n-1} + a_0x^n$. Prove q is irreducible in $F[x]$.
66. Let $p \in \mathbb{N}^+$ be a prime integer, and let $g = x^{p-1} + x^{p-2} + \dots + x^2 + x + 1$. This problem will prove that g is irreducible in $\mathbb{Q}[x]$. Recall (Exercise 20) that the map $\phi : \mathbb{Q}[x] \rightarrow \mathbb{Q}[x]$, given by $\phi(f) = f(x+1)$ for $f \in \mathbb{Q}[x]$, is a ring isomorphism. (a) Compute the coefficients of the polynomial $g^* = g(x+1)$. [Hint: To start, apply ϕ to the equation $x^p - 1 = (x-1)g$.] (b) Explain why g is irreducible in $\mathbb{Q}[x]$ iff g^* is irreducible in $\mathbb{Q}[x]$. (c) Use a theorem from the text to prove g^* is irreducible in $\mathbb{Q}[x]$.
67. Does Eisenstein's criterion hold if we only assume that p is a positive integer (not necessarily prime)? Prove or give a justified counterexample. What happens if p is a product of two or more distinct primes?
68. (a) Find $a, b, c, d \in \mathbb{Q}$ such that $f = a + bx + cx^2 + dx^3 \in \mathbb{Q}[x]$ satisfies $f(0) = -4$, $f(1) = -2$, $f(2) = 6$, and $f(1/2) = -3$. (b) Find $a, b, c \in \mathbb{Z}_7$ such that $g = a + bx + cx^2 \in \mathbb{Z}_7[x]$ satisfies $g(1) = 5$, $g(2) = 3$, and $g(3) = 0$. (c) Find $a, b, c \in \mathbb{Z}_{13}$ such that $h = a + bx + cx^2 \in \mathbb{Z}_{13}[x]$ satisfies $h(2) = 12$, $h(4) = 10$, and $h(6) = 9$.
69. *Secret Sharing.* A set of n people wish to share information about a master secret

in such a way that any k of the n people can pool their knowledge to recover the secret, but no subset of fewer than k people can gain any knowledge about the secret by working together. This goal can be achieved using polynomials, as follows. Pick a prime $p > n$, and encode the master secret as a value $a_0 \in \mathbb{Z}_p$. Choose random $a_1, \dots, a_{k-1} \in \mathbb{Z}_p$, and let $g = \sum_{i=0}^{k-1} a_i x^i \in \mathbb{Z}_p[x]$. Number the people 1 to n arbitrarily, and give person j the pair $(j, g(j))$. (a) Show how any subset of k or more people can use their collective knowledge to recover a_0 . (b) Prove that the collective knowledge of any subset of $k - 1$ or fewer people gives no information about the secret.

70. *Chinese Remainder Theorem.* There is an analogue of the Lagrange interpolation formula that lets us solve systems of congruences with pairwise relatively prime moduli. Assume n_1, \dots, n_k are positive integers with $\gcd(n_i, n_j) = 1$ for all $i \neq j$. Set $N = n_1 n_2 \cdots n_k$. (a) Prove: for every list (b_1, \dots, b_k) with each $b_i \in \mathbb{Z}_{n_i}$, there exists a unique $x \in \mathbb{Z}_N$ such that $x \bmod n_i = b_i$ for $1 \leq i \leq k$. [Hint: To prove existence, define $q_i = \prod_{j \neq i} n_j$ for $1 \leq i \leq k$. Note $\gcd(q_i, n_i) = 1$, so Exercise 32 gives an integer $r_i \in \mathbb{Z}_{n_i}$ with $(q_i r_i) \bmod n_i = 1$. Define $p_i = q_i r_i \in \mathbb{Z}_N$; what is $p_i \bmod n_i$ and $p_i \bmod n_j$ for $j \neq i$? Use the p_i 's and b_i 's to construct x .] (b) Use the method sketched in (a) to find $x \in \mathbb{Z}_{1001}$ solving $x \bmod 7 = 5$, $x \bmod 11 = 4$, and $x \bmod 13 = 1$.
71. For each polynomial, use Kronecker's algorithm to factor it in $\mathbb{Q}[x]$ or prove that it is irreducible in $\mathbb{Q}[x]$. (a) $x^4 - x^2 + 2x - 1$; (b) $x^5 - 4x + 2$; (c) $x^5 - 3x^4 + 10x^3 - 10x^2 + 24x - 7$.
72. (a) Prove that $(1, \sqrt{2}, \sqrt{3})$ is a \mathbb{Q} -linearly independent list. (b) Deduce that $(1, \sqrt{2}, \sqrt{3}, \sqrt{6})$ is a \mathbb{Q} -linearly independent list. (c) Let $z = \sqrt{2} + \sqrt{3}$. Prove, as asserted in the text, that $(1, z, z^2, z^3)$ is \mathbb{Q} -linearly independent.
73. Viewing \mathbb{C} as a \mathbb{Q} -algebra, find (with proof) the minimal polynomials over \mathbb{Q} of each of these complex numbers: (a) $-i$; (b) $\sqrt[3]{7}$; (c) $e^{\pi i/4}$; (d) $\sqrt{5} + 3i$; (e) $\sqrt{2} + \sqrt[3]{5}$.
74. Find the minimal polynomial of each matrix over \mathbb{R} (where $a, b, c, d \in \mathbb{R}$):
 (a) $\begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$; (b) $\begin{bmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{bmatrix}$; (c) $\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$; (d) $\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -d & -c & -b \end{bmatrix}$;
 (e) $\begin{bmatrix} a & b \\ 0 & d \end{bmatrix}$.
75. Let F be a field. (a) Show that a polynomial z in the F -algebra $F[x]$ is algebraic over F iff z is a constant. (b) Does (a) hold for all z in the F -algebra $F[[x]]$?
76. *Formal Derivatives.* Given a commutative ring R and $f = \sum_{n \geq 0} f_n x^n \in R[[x]]$, define the *formal derivative* of f by setting $D(f) = \sum_{n \geq 1} n f_n x^{n-1} = ((n+1)f_{n+1} : n \in \mathbb{N})$, where $n f_n$ denotes the sum of n copies of f_n in R . (a) Prove that D is R -linear: for all $f, g \in R[[x]]$ and all $c \in R$, $D(f+g) = D(f) + D(g)$ and $D(cf) = cD(f)$. (b) Prove the *Leibniz product rule* for D : for all $f, g \in R[[x]]$, $D(fg) = D(f)g + fD(g)$. (c) Prove the *formal chain rule* for D : for all $f \in R[x]$ and $g \in R[[x]]$, $D(f(g)) = (D(f))(g) \cdot D(g)$. Here $f(g)$ is the evaluation of the polynomial f at $x = g$, and similarly for $(D(f))(g)$. (d) Which results in (a), (b), (c) hold for non-commutative R ?
77. *Repeated Roots.* Let F be a field, $p \in F[x]$, and $c \in F$. We say c is a *repeated root* of p iff $(x - c)^2 | p$ in $F[x]$. (a) Prove c is a repeated root of p iff $p(c) = 0 = D(p)(c)$ (see the previous exercise for the definition of $D(p)$). (b) Prove: if p has a repeated root in some field K containing F as a subfield, then $\gcd(p, D(p))$ (computed in

$F[x]$) is not 1. (The converse of (b) is also true, but you will need Exercise 34 in Chapter 12 to prove it.)

78. True or false? Explain each answer. (a) For all commutative rings R and all $g \in R[x]$ of degree $n > 0$, g has at most n roots in R . (b) For all fields F and all $p \in F[x]$, if p is irreducible in $F[x]$ then p has no root in F . (c) For all fields F and all $p \in F[x]$, if p has no root in F then p is irreducible in $F[x]$. (d) For all fields F and all $p \in F[x]$, if $p(c) = 0$ for all $c \in F$, then $p = 0$. (e) For all irreducible $f \in \mathbb{Q}[x]$, the only divisors of f in $\mathbb{Q}[x]$ are 1 and f . (f) The polynomial $x^5 + x^4 + x^3 + x^2 + x + 1$ is irreducible in $\mathbb{Q}[x]$. (g) For all $b, c \in \mathbb{Z}$, every rational root of $2x^3 + bx^2 + cx + 3$ must lie in the set $\{\pm 1, \pm 2, \pm 1/3, \pm 2/3\}$. (h) For all fields F and nonzero $f, g \in F[x]$, $\gcd(f, g) = 1$ iff $1 = uf + vg$ for some $u, v \in F[x]$. (i) For all rings R and all nonzero $f, g \in R[x]$, $\deg(fg) = \deg(f) + \deg(g)$. (j) For all $f, g \in \mathbb{Z}_2[x]$, if $f|g$ and $g|f$ then $f = g$. (k) For all even $a \in \mathbb{Z}$ and odd $b, c \in \mathbb{Z}$, $x^3 + ax^2 + bx + c$ is irreducible in $\mathbb{Q}[x]$. (l) Every monic $f \in \mathbb{Q}[x]$ with $\deg(f) \geq 1$ can be factored as $f = p_1 p_2 \cdots p_k$, where each p_i is irreducible in $\mathbb{Q}[x]$, and the only other factorizations of f into irreducible factors are obtained by reordering the p_i 's. (m) For all $p, q \in \mathbb{Q}[x]$, if $p(m) = q(m)$ for all $m \in \mathbb{Z}$, then $p = q$. (n) The minimal polynomial of an algebraic element in a commutative F -algebra must be irreducible in $F[x]$. (o) For all fields F and all monic $d, f, g \in F[x]$, $\gcd(f, g) = d$ iff $d = uf + vg$ for some $u, v \in F[x]$.

This page intentionally left blank

Part II

Matrices

This page intentionally left blank

Basic Matrix Operations

Our aim in this chapter is to treat some elementary material about matrices from a somewhat more advanced viewpoint. We begin by giving formal definitions of matrices, ordered n -tuples, row vectors, and column vectors. We then give several equivalent formulations of the familiar matrix operations (addition, scalar multiplication, matrix multiplication, transpose, etc.) and derive their basic properties. We also discuss how to implement elementary row and column operations by multiplying a matrix on the left or right by certain “elementary matrices.”

Throughout this chapter, the letter F will denote an arbitrary *field*; elements of F will be called *scalars*. However, most of our results on matrices extend readily to the case where F is a commutative ring, or even to the case of arbitrary rings F . In this more general setting, vector spaces are replaced by modules (see Chapter 17), and nonzero scalars in the field F are replaced by units (invertible elements) in the ring F . For those readers who are interested in this level of generality, we include annotations in square brackets at any places in the chapter that explicitly require the assumptions that F is a field, or that F is commutative.

4.1 Formal Definition of Matrices and Vectors

For each positive integer n , write $[n]$ for the set $\{1, 2, \dots, n\}$. Informally, an n -tuple with entries in F is an ordered list $x = (x_1, \dots, x_n)$ of scalars in F . Formally, we regard this n -tuple as the function $x : [n] \rightarrow F$ given by $x(i) = x_i$ for $i \in [n]$. We will freely use either notation (lists or functions) to describe n -tuples. Let F^n be the set of all n -tuples with entries in F .

Informally, an $m \times n$ matrix A over F is an array of scalars consisting of m rows and n columns. The element of F in row i and column j of A is called the i, j -entry of A and is denoted $A_{i,j}$ or A_{ij} or a_{ij} . Formally, let $[m] \times [n]$ be the set of ordered pairs $\{(i, j) : i \in [m], j \in [n]\}$. We define an $m \times n$ matrix to be a function $A : [m] \times [n] \rightarrow F$, where $A(i, j)$ is the i, j -entry of A . Let $M_{m,n}(F)$ be the set of all $m \times n$ matrices over F , and let $M_n(F)$ be the set of all $n \times n$ (square) matrices over F .

A *column vector* of length n is an $n \times 1$ matrix. A *row vector* of length n is a $1 \times n$ matrix. Formally, a column vector is a function $x_c : [n] \times [1] \rightarrow F$, whereas a row vector is a function $x_r : [1] \times [n] \rightarrow F$. There is a natural bijection (one-to-one correspondence) between column vectors and n -tuples that sends x_c to the n -tuple $(x_c(1, 1), x_c(2, 1), \dots, x_c(n, 1))$. Similarly, there is a bijection between row vectors and n -tuples that sends x_r to the n -tuple $(x_r(1, 1), x_r(1, 2), \dots, x_r(1, n))$. Using these bijections, we can regard column vectors and row vectors as n -tuples whenever it is convenient to do so.

Given a matrix $A \in M_{m,n}(F)$ and $j \in [n]$, we write $A^{[j]}$ to denote the j 'th column of A , which can also be regarded as the m -tuple $(A(1, j), A(2, j), \dots, A(m, j))$. For $i \in [m]$, we write $A_{[i]}$ to denote the i 'th row of A , which can be regarded as the n -tuple

$(A(i, 1), A(i, 2), \dots, A(i, n))$. Any matrix A is completely determined by the ordered list of its columns $(A^{[1]}, \dots, A^{[n]})$ and also by the ordered list of its rows $(A_{[1]}, \dots, A_{[m]})$. Formally, this remark means that we have a bijection from $M_{m,n}(F)$ to $(F^m)^n$ and another bijection between $M_{m,n}(F)$ and $(F^n)^m$.

Now we define some special matrices. The rectangular *zero matrices* $0_{m,n} \in M_{m,n}(F)$ are defined by $0_{m,n}(i, j) = 0$ for $1 \leq i \leq m$ and $1 \leq j \leq n$. We write 0_n for $0_{n,n}$. The square *identity matrices* $I_n \in M_n(F)$ are defined by $I_n(i, i) = 1$ and $I_n(i, j) = 0$ for $i \neq j$ (where $1 \leq i, j \leq n$). The *unit matrices* $J_{m,n} \in M_{m,n}(F)$ are defined by $J_{m,n}(i, j) = 1$ for $1 \leq i \leq m$ and $1 \leq j \leq n$. We omit the subscripts on these matrices whenever they are understood from context.

Recall that two functions f and g are *equal* iff they have the same domain, the same codomain, and $f(x) = g(x)$ for all x in the common domain. Applying this remark to the formal definition of matrices, we see that an $m \times n$ matrix A equals an $m' \times n'$ matrix A' iff $m = m'$ and $n = n'$ and $A(i, j) = A'(i, j)$ for all $1 \leq i \leq m$ and $1 \leq j \leq n$. Similarly, $v \in F^n$ equals $v' \in F^{n'}$ iff $n = n'$ and $v(i) = v'(i)$ for $1 \leq i \leq n$.

4.2 Vector Spaces of Functions

We are about to give “entry-by-entry” definitions of algebraic operations on matrices and vectors. Before doing so, it is helpful to consider a more general situation. Let F be the field of scalars, let S be an arbitrary set, and let V be the set of all functions $g : S \rightarrow F$. By introducing “pointwise operations” on functions, we can turn the set V into an F -vector space. [When F is only a ring, V will be a left F -module.] First, we define an addition operation $+ : V \times V \rightarrow V$. Let $g, h : S \rightarrow F$ be two elements of V . Define $g + h : S \rightarrow F$ by setting $(g + h)(x) = g(x) + h(x)$ for all $x \in S$. Note that the plus symbol on the right side denotes addition in the given field F , while the plus symbol on the left side is the new addition of functions that is being defined. Since F is closed under addition, $g(x) + h(x)$ always lies in F , so that $g + h$ is a function from S to F . In other words, we have the closure condition: for all $g, h \in V$, $g + h \in V$. Similarly, the other additive axioms in the definition of a vector space [or module] follow from the corresponding axioms for the field [or ring] F . For instance, to confirm that $(g + h) + k = g + (h + k)$ for all $g, h, k \in V$, we check that both sides (which are functions) agree at every $x \in S$:

$$\begin{aligned} [(g + h) + k](x) &= (g + h)(x) + k(x) = (g(x) + h(x)) + k(x) \\ &= g(x) + (h(x) + k(x)) = g(x) + (h + k)(x) = [g + (h + k)](x). \end{aligned} \quad (4.1)$$

Note that this calculation used associativity of addition in F . A similar computation shows that $g + h = h + g$ for all $g, h \in V$. One may check that the additive identity of V is the *zero function* $0_V : S \rightarrow F$ given by $0_V(x) = 0_F$ for all $x \in S$. The additive inverse of $f \in V$ is the function $(-f) : S \rightarrow F$ defined by $(-f)(x) = -(f(x))$, where the minus on the right side denotes the additive inverse of $f(x)$ in F .

Next, we define a scalar multiplication operation $\cdot : F \times V \rightarrow V$. Given any $g : S \rightarrow F$ in V and any scalar $c \in F$, define the function $c \cdot g : S \rightarrow F$ by the equation $(c \cdot g)(x) = c \cdot (g(x))$ for $x \in S$; here, the \cdot on the right side is multiplication in the given field [or ring] F . The remaining axioms for a vector space [or module] can now be routinely checked, using the corresponding properties of the field [or ring] F . We can now consider finite linear combinations of elements of V , as in any vector space [or module]. Explicitly, if $f_i \in V$ and

$c_i \in F$, then $\sum_{i=1}^n c_i f_i \in V$ is the function from S to F sending x to $\sum_{i=1}^n c_i f_i(x) \in F$ for all $x \in S$.

From now on, we assume S is a *finite* set. For each $x \in S$, we can define a function $e_x : S \rightarrow F$ by letting $e_x(x) = 1_F$ and $e_x(y) = 0_F$ for all $y \in S$ with $y \neq x$. Each e_x is an element of V , so the set $X = \{e_x : x \in S\}$ is a subset of V . There is a bijection from S to X given by $x \mapsto e_x$ (using the fact that $1_F \neq 0_F$ in F [which also holds if F is a nonzero ring]).

We now show that X is a *basis* for the F -vector space V . To prove linear independence, assume $\sum_{x \in S} c_x e_x = 0_V$ where each $c_x \in F$. Both sides of this identity are functions with domain S . Evaluate each side at some fixed $x_0 \in S$. Using the fact that $e_x(x_0) = 0$ for $x \neq x_0$ and $e_{x_0}(x_0) = 1_F$, we obtain $c_{x_0} = 0_V(x_0) = 0_F$. This holds for every $x_0 \in S$, so linear independence is proved. To show that X spans V , let $f : S \rightarrow F$ be any element of V . The spanning assertion will follow from the identity $\sum_{x \in X} f(x)e_x = f$, which exhibits f as an F -linear combination of elements of X . To prove this equality between two functions, it suffices to show that both sides agree at any $x_0 \in S$. Reasoning as above, we see that

$$\left(\sum_{x \in X} f(x)e_x \right)(x_0) = \sum_{x \in X} f(x)e_x(x_0) = f(x_0),$$

completing the proof. As a corollary, observe that $\dim(V) = |X| = |S|$.

4.3 Matrix Operations via Entries

We now define the vector space operations on matrices and n -tuples as special cases of the constructions in the previous subsection. First consider the set F^n of n -tuples of scalars from F . Formally, F^n is the set of all functions from the set $S = [n]$ to F . Therefore, F^n is an F -vector space via the pointwise operations on functions introduced in the last section. Reverting to the informal “list” notation for n -tuples, the definitions of these operations read:

$$(x_1, x_2, \dots, x_n) + (y_1, y_2, \dots, y_n) = (x_1 + y_1, x_2 + y_2, \dots, x_n + y_n) \quad \text{for } x_i, y_i \in F;$$

$$c(x_1, x_2, \dots, x_n) = (cx_1, cx_2, \dots, cx_n) \quad \text{for } c, x_i \in F.$$

The basis vector e_i (where $1 \leq i \leq n$) is the function such that $e_i(i) = 1_F$ and $e_i(j) = 0_F$ for $j \neq i$. Using the list notation, e_i is a list that has 1_F in position i and 0_F elsewhere: $e_i = (0, \dots, 1, \dots, 0)$. Our general results show that $X = (e_1, \dots, e_n)$ is an ordered basis for F^n ; we call this the *standard ordered basis for F^n* . In the list notation, we have

$$(x_1, \dots, x_n) = x_1 e_1 + \cdots + x_n e_n.$$

Note that the vector e_i depends on n (the domain of e_i is $[n]$), and this dependence is not indicated in the notation e_i . In practice, n is understood from context, so no ambiguity occurs. We have $\dim(F^n) = n$.

Now consider the set $M_{m,n}(F)$ of $m \times n$ matrices over F . Letting $S = [m] \times [n]$, the formal definition of matrices shows that $M_{m,n}(F)$ is the set of all functions from S to F . We know this is a vector space via pointwise operations on functions. In terms of the entries of the matrices, the definitions of the operations are:

$$(A + B)(i, j) = A(i, j) + B(i, j), \quad (cA)(i, j) = c(A(i, j))$$

$$(A, B \in M_{m,n}(F), c \in F, 1 \leq i \leq m, 1 \leq j \leq n). \quad (4.2)$$

The basis vector $e_{(i,j)} = e_{ij}$ is the function from $[m] \times [n]$ to F such that $e_{ij}(i, j) = 1_F$ and $e_{ij}(i', j') = 0_F$ if $i' \neq i$ or $j' \neq j$. Using the array notation for matrices, e_{ij} is the array that has 1_F in row i and column j , and zeroes in all other positions. We know that the set $X = \{e_{ij} : 1 \leq i \leq m, 1 \leq j \leq n\}$ is a basis for $M_{m,n}(F)$, which is called the *standard basis for $M_{m,n}(F)$* . If A is a matrix with i, j -entry a_{ij} , we have the matrix identity $A = \sum_{i=1}^m \sum_{j=1}^n a_{ij} e_{ij}$. For example,

$$\begin{bmatrix} 2 & -1 \\ 0 & 6 \end{bmatrix} = 2 \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} + (-1) \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} + 0 \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} + 6 \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

As above, note that the meaning of e_{ij} does depend on m and n . We have $\dim(M_{m,n}(F)) = |X| = mn$.

For each $n \in \mathbb{N}^+$, we now have the three n -dimensional vector spaces $M_{n,1}(F)$, $M_{1,n}(F)$, and F^n , which respectively consist of column vectors of length n , row vectors of length n , and n -tuples. The natural bijections between these sets are readily seen to be vector space isomorphisms. For $i \in [n]$, the basis vectors $e_{i1} \in M_{n,1}(F)$, $e_{1i} \in M_{1,n}(F)$, and $e_i \in F^n$ correspond to each other under these isomorphisms.

We now define the operations of *matrix multiplication*, *matrix-vector multiplication*, *matrix transpose*, and *conjugate-transpose* entry by entry. Let $A \in M_{m,n}(F)$ and $B \in M_{n,p}(F)$. The *matrix product* AB is the $m \times p$ matrix defined by

$$(AB)(i, j) = \sum_{k=1}^n A(i, k)B(k, j) \quad (1 \leq i \leq m, 1 \leq j \leq p).$$

By letting $m = 1$ or $p = 1$, we recover formulas for multiplying a matrix on the left or right by a row vector or column vector (respectively). Specifically, if $A \in M_{m,n}(F)$ and $v \in F^n = M_{n,1}(F)$, then $Av \in F^m = M_{m,1}(F)$ is defined by

$$(Av)(i) = \sum_{k=1}^n A(i, k)v(k) \text{ for } i \in [m].$$

(We are using the identification conventions mentioned earlier, so that $(Av)(i) = (Av)(i, 1)$ and $v(k) = v(k, 1)$ here.) If, instead, $w \in F^m = M_{1,m}(F)$ and $A \in M_{m,n}(F)$, then $wA \in F^n = M_{1,n}(F)$ is defined by

$$(wA)(j) = \sum_{k=1}^m w(k)A(k, j) \text{ for } j \in [n].$$

For example, if $A = \begin{bmatrix} 2 & -1 & 4 \\ 0 & 1 & 3 \end{bmatrix} \in M_{2,3}(\mathbb{R})$, $B = \begin{bmatrix} 3 & 2 \\ 1 & 1 \\ -2 & 5 \end{bmatrix} \in M_{3,2}(\mathbb{R})$, and $v = (4, 0, 1) \in \mathbb{R}^3$, then

$$AB = \begin{bmatrix} -3 & 23 \\ -5 & 16 \end{bmatrix}, \quad BA = \begin{bmatrix} 6 & -1 & 18 \\ 2 & 0 & 7 \\ -4 & 7 & 7 \end{bmatrix}, \quad vB = [10 \ 13], \quad Av = \begin{bmatrix} 12 \\ 3 \end{bmatrix}.$$

Note that we view the 3-tuple v as a row vector when computing vB , but we view v as a column vector when computing Av .

For general $A \in M_{m,n}(F)$, the *transpose* of A , denoted A^T , is the $n \times m$ matrix defined by

$$A^T(i, j) = A(j, i) \quad (1 \leq i \leq n, 1 \leq j \leq m).$$

For $F = \mathbb{C}$, the *conjugate-transpose* of A , denoted A^* , is the $n \times m$ complex matrix defined by

$$A^*(i, j) = \overline{A(j, i)} \quad (1 \leq i \leq n, 1 \leq j \leq m).$$

Here the bar denotes complex conjugation: $\overline{x+iy} = x-iy$ for $x, y \in \mathbb{R}$. For example, if $A = \begin{bmatrix} 2+i & -3i & 4 \\ 1-2i & 0 & e^{\pi i/4} \end{bmatrix} \in M_{2,3}(\mathbb{C})$, then

$$A^T = \begin{bmatrix} 2+i & 1-2i & 4 \\ -3i & 0 & e^{\pi i/4} \\ 4 & e^{\pi i/4} & \end{bmatrix}, \quad A^* = \begin{bmatrix} 2-i & 1+2i & 0 \\ 3i & 0 & \\ 4 & e^{-\pi i/4} & \end{bmatrix}.$$

4.4 Properties of Matrix Multiplication

We now use the definition of the product of two matrices to derive some of the fundamental algebraic properties of matrix multiplication. Let $A, A' \in M_{m,n}(F)$ and $B, B' \in M_{n,p}(F)$. First, we have the two *distributive laws* $A(B+B') = AB+AB'$ and $(A+A')B = AB+A'B$. To prove the first of these laws, compute the i, j -entry of each side, for $i \in [m]$ and $j \in [p]$:

$$\begin{aligned} [A(B+B')](i, j) &= \sum_{k=1}^n A(i, k)(B+B')(k, j) = \sum_{k=1}^n A(i, k)[B(k, j) + B'(k, j)] \\ &= \sum_{k=1}^n [A(i, k)B(k, j) + A(i, k)B'(k, j)] \\ &= \sum_{k=1}^n A(i, k)B(k, j) + \sum_{k=1}^n A(i, k)B'(k, j) \\ &= (AB)(i, j) + (AB')(i, j) = (AB + AB')(i, j). \end{aligned}$$

The second distributive law is proved similarly. Next, for $A \in M_{m,n}(F)$ and $B \in M_{n,p}(F)$ and $c \in F$ [and F commutative], we have $c(AB) = (cA)B = A(cB)$. This is proved by showing that the i, j -entry of all these expressions is $\sum_{k=1}^n cA(i, k)B(k, j)$. We also have $AI_n = A = I_m A$. For instance, the second equality follows from

$$(I_m A)(i, j) = \sum_{k=1}^m I_m(i, k)A(k, j) = A(i, j) \text{ for } i \in [m] \text{ and } j \in [n],$$

since $I_m(i, k)$ is 1 for $k = i$ and 0 otherwise.

Finally, letting $A \in M_{m,n}(F)$, $B \in M_{n,p}(F)$, and $C \in M_{p,q}(F)$, we have the *associative law* $(AB)C = A(BC) \in M_{m,q}(F)$. To prove associativity, choose any $i \in [m]$ and $j \in [q]$. On one hand,

$$\begin{aligned} [(AB)C](i, j) &= \sum_{k'=1}^p (AB)(i, k')C(k', j) = \sum_{k'=1}^p \left(\sum_{k=1}^n A(i, k)B(k, k') \right) C(k', j) \\ &= \sum_{k'=1}^p \sum_{k=1}^n (A(i, k)B(k, k'))C(k', j). \end{aligned}$$

On the other hand,

$$\begin{aligned}[A(BC)](i,j) &= \sum_{k=1}^n A(i,k)(BC)(k,j) = \sum_{k=1}^n A(i,k) \left(\sum_{k'=1}^p B(k,k')C(k',j) \right) \\ &= \sum_{k=1}^n \sum_{k'=1}^p A(i,k)(B(k,k')C(k',j)).\end{aligned}$$

The two final expressions are equal, since the finite sums (indexed by k and k') can be interchanged and since multiplication is associative in the field F . [The proof holds for any ring F , possibly non-commutative.]

Letting $m = n = p = q$, the results proved above show that the F -vector space $M_n(F)$ is an associative F -algebra with identity (see §1.3). This means that $M_n(F)$ is an F -vector space and a ring (with identity I_n) such that $c(AB) = (cA)B = A(cB)$ for all $c \in F$ and all $A, B \in M_n(F)$.

We conclude this section with some properties of the transpose and conjugate-transpose operations. Let $A, B \in M_{m,n}(F)$ and $C \in M_{n,p}(F)$. We have $(A + B)^T = A^T + B^T$ since the i, j -entry of both sides is $A(j, i) + B(j, i)$ for all $i \in [n]$ and $j \in [m]$. For $c \in F$, we have $(cA)^T = c(A^T)$ since the i, j -entry of both sides is $cA(j, i)$. Finally, we have $(AC)^T = C^T A^T$ since, for all $i \in [p]$ and $j \in [m]$,

$$(AC)^T(i, j) = (AC)(j, i) = \sum_{k=1}^n A(j, k)C(k, i) = \sum_{k=1}^n C^T(i, k)A^T(k, j) = (C^T A^T)(i, j).$$

[Commutativity of multiplication in F was used in this proof.] The conjugate-transpose operation has similar properties: $(A + B)^* = A^* + B^*$; $(cA)^* = \bar{c}(A^*)$; and $(AC)^* = C^* A^*$. The proofs are similar to those just given; for instance, the second property is true because, for $i \in [n]$ and $j \in [m]$,

$$(cA)^*(i, j) = \overline{(cA)(j, i)} = \overline{c(A(j, i))} = \bar{c} \cdot \overline{A(j, i)} = \bar{c}(A^*(i, j)) = (\bar{c}A^*)(i, j).$$

We also have $(A^T)^T = A$ and, for $F = \mathbb{C}$, $(A^*)^* = A$.

4.5 Generalized Associativity

We now prove a general associativity result that applies to products of three or more matrices whose dimensions match. Informally, the result says that no parentheses are required to evaluate such a product unambiguously (although the order of the factors certainly matters in general). Formally, suppose we are given $s \geq 1$ matrices A_1, \dots, A_s and integers n_0, \dots, n_s such that $A_i \in M_{n_{i-1}, n_i}(F)$ for each i . Suppose also that we are given any complete parenthesization of the sequence $A_1 A_2 \cdots A_s$, which specifies exactly how these matrices are to be combined via the binary operation of matrix multiplication. For example, if $s = 4$, there are five possible complete parenthesizations:

$$A_1(A_2(A_3A_4)), ((A_1A_2)A_3)A_4, (A_1A_2)(A_3A_4), A_1((A_2A_3)A_4), (A_1(A_2A_3))A_4.$$

We define the *standard* complete parenthesization to be the $n_0 \times n_s$ matrix

$$\prod_{i=1}^s A_i = (\cdots (((A_1A_2)A_3)A_4) \cdots)A_s.$$

More formally, we give the recursive definition $\prod_{i=1}^1 A_i = A_1$ and $\prod_{i=1}^s A_i = (\prod_{i=1}^{s-1} A_i)A_s$ for all $s > 1$. We will show that all complete parenthesizations of the given sequence evaluate to $\prod_{i=1}^s A_i$. We prove the result by induction on s . The cases $s \leq 2$ are immediate, while the case $s = 3$ is the ordinary associative law proved in the previous section. Now assume $s > 3$ and the result is known to hold for smaller values of s . Given any complete parenthesization of $A_1 \cdots A_s$, there exists a unique index $t < s$ such that the *last* binary product operation involved in the computation involves multiplying some complete parenthesization of $A_1 \cdots A_t$ by some complete parenthesization of $A_{t+1} \cdots A_s$. For example, for the five parenthesizations listed above in the case $s = 4$, we have $t = 1, t = 3, t = 2, t = 1$, and $t = 3$, respectively. By induction, the complete parenthesization of the first t factors evaluates to the $n_0 \times n_t$ matrix $B = \prod_{i=1}^t A_i$, while the complete parenthesization of the last $s - t$ factors evaluates to the $n_t \times n_s$ matrix $C = \prod_{j=t+1}^s A_j$. If $t = s - 1$, this second product is just A_s , and then the given complete parenthesization of $A_1 \cdots A_s$ evaluates to $BA_s = \prod_{i=1}^s A_i$ by definition. Otherwise, in the case $t < s - 1$, set $D = \prod_{j=t+1}^{s-1} A_j$. We have $C = DA_s$ by definition, so that the given complete parenthesization of $A_1 \cdots A_s$ evaluates to $BC = B(DA_s)$. By associativity for three factors and induction, this equals $(BD)A_s = (\prod_{i=1}^{s-1} A_i)A_s = \prod_{i=1}^s A_i$. This completes the proof.

There is a formula for the entries of $\prod_{u=1}^s A_u$ that extends the formula for a product of two matrices. Specifically, the i, j entry of $\prod_{u=1}^s A_u$ is

$$\sum_{k_1=1}^{n_1} \sum_{k_2=1}^{n_2} \cdots \sum_{k_{s-1}=1}^{n_{s-1}} A_1(i, k_1) A_2(k_1, k_2) A_3(k_2, k_3) \cdots A_{s-1}(k_{s-2}, k_{s-1}) A_s(k_{s-1}, j).$$

This can be proved by induction on s , where the base case $s = 2$ is precisely the definition of matrix multiplication. For $s > 2$, this definition shows that the i, j -entry of $\prod_{u=1}^s A_u = (\prod_{u=1}^{s-1} A_u)A_s$ is

$$\sum_{k=1}^{n_{s-1}} \left[\prod_{u=1}^{s-1} A_u \right] (i, k) A_s(k, j),$$

which, by induction, equals

$$\sum_{k=1}^{n_{s-1}} \left(\sum_{k_1=1}^{n_1} \cdots \sum_{k_{s-2}=1}^{n_{s-2}} A_1(i, k_1) A_2(k_1, k_2) \cdots A_{s-1}(k_{s-2}, k) \right) A_s(k, j).$$

Renaming the summation variable k to be k_{s-1} , using the distributive law in F , and reordering the finite summations, we obtain the stated formula for the i, j -entry of a product of s factors.

Now that we know the products are unambiguously defined, one may prove that $(A_1 A_2 \cdots A_s)^T = A_s^T \cdots A_2^T A_1^T$ and, when $F = \mathbb{C}$, $(A_1 A_2 \cdots A_s)^* = A_s^* \cdots A_2^* A_1^*$. We already proved these identities for $s = 2$, and the general case follows by induction using generalized associativity. With generalized associativity in hand, we can also define the *positive powers* of a square matrix $A \in M_n(F)$. For each integer $s \geq 1$, we let $A^s = \prod_{i=1}^s A_i$ where every $A_i = A$. Informally, A^s is the product of s copies of A . We also define $A^0 = I_n$.

4.6 Invertible Matrices

A matrix A with entries in F is *invertible* iff A is square (say $n \times n$) and there exists a matrix $B \in M_n(F)$ with $AB = I_n = BA$. B is called an *inverse* of A . If such a matrix B

exists, it is unique. For if $B_1 \in M_n(F)$ also satisfies $AB_1 = I_n = B_1A$, then associativity of matrix multiplication gives

$$B_1 = B_1I_n = B_1(AB) = (B_1A)B = I_nB = B.$$

For an invertible matrix A , we write A^{-1} to denote the unique inverse of A .

For all $n \geq 1$ and all matrices $C, D \in M_n(F)$, we can check that C^{-1} exists and that $C^{-1} = D$ by merely verifying the defining condition $CD = I_n = DC$. For example, let us show that *if $A \in M_n(F)$ is invertible, then A^{-1} is also invertible, and $(A^{-1})^{-1} = A$* . By definition, $A^{-1}A = I_n = AA^{-1}$. By the remark at the beginning of this paragraph (taking $C = A^{-1}$ and $D = A$), we obtain the claimed result on the invertibility of A^{-1} . Similarly, I_n is invertible with $I_n^{-1} = I_n$. To see this, take $C = D = I_n$ above and note that $I_nI_n = I_n = I_nI_n$.

For another example, suppose U and V are invertible matrices in $M_n(F)$. We claim the product UV is invertible also, with inverse $(UV)^{-1} = V^{-1}U^{-1}$. To see this, compute

$$(UV)(V^{-1}U^{-1}) = U(VV^{-1})U^{-1} = UI_nU^{-1} = UU^{-1} = I_n$$

and, similarly, $(V^{-1}U^{-1})(UV) = I_n$. The remark above, with $C = UV$ and $D = V^{-1}U^{-1}$, completes the proof. More generally, given k invertible matrices $U_1, \dots, U_k \in M_n(F)$, the product $U_1U_2 \cdots U_k$ is also invertible, and $(U_1U_2 \cdots U_k)^{-1} = U_k^{-1} \cdots U_2^{-1}U_1^{-1}$. One can prove this by induction on k (the case $k = 2$ was done above). Note that generalized associativity was used heavily throughout this paragraph. If we take all U_i 's equal to a given invertible matrix A , we see that $(A^k)^{-1} = (A^{-1})^k$ for all $k \geq 1$. For invertible A , we define the *negative power* A^{-k} to be $(A^{-1})^k = (A^k)^{-1}$.

For all fields F and all positive integers n , let $\text{GL}_n(F)$ be the set of all invertible matrices in $M_n(F)$. The set $\text{GL}_n(F)$ is a group under matrix multiplication, which is almost always non-commutative. The group axioms (closure, associativity, identity, and inverses) follow from the results proved above. $\text{GL}_n(F)$ is called the *general linear group of degree n*.

For any matrix $A \in M_n(F)$, A is invertible iff the transpose A^T is invertible, in which case $(A^T)^{-1} = (A^{-1})^T$. To see this, first assume A^{-1} exists. Then

$$(A^T)(A^{-1})^T = (A^{-1}A)^T = (I_n)^T = I_n = (I_n)^T = (AA^{-1})^T = (A^{-1})^T(A^T),$$

so the conclusion follows by taking $C = A^T$ and $D = (A^{-1})^T$ in the remark above. So the forward implication holds for *all* matrices A . Applying the forward implication to the matrix A^T , we conclude that the invertibility of A^T implies the invertibility of $(A^T)^T = A$, which gives the reverse implication for the original A . An entirely analogous argument, using the fact that $(I_n)^* = I_n$, shows that *for all $A \in M_n(\mathbb{C})$, A is invertible iff A^* is invertible, in which case $(A^*)^{-1} = (A^{-1})^*$* .

[The rest of this section applies to fields and commutative rings F .] The next few paragraphs assume familiarity with determinants, which will be studied in detail in Chapter 5. We recall that for every square matrix $A \in M_n(F)$, there is an associated scalar $\det(A) \in F$ called the *determinant of A*. One of the most well-known theorems of linear algebra states that *for any field F and any $A \in M_n(F)$, A^{-1} exists iff $\det(A) \neq 0_F$* . In §5.11, we will prove the more general fact that for any commutative ring R and any $A \in M_n(R)$, A^{-1} exists in $M_n(R)$ iff $\det(A)$ is an invertible element of the ring R . The proof will provide explicit formulas (involving determinants) for every entry of A^{-1} , when the inverse exists. We will also prove in §5.13 that for any commutative ring R and all $A, B \in M_n(R)$, $\det(AB) = \det(A)\det(B)$.

In our definition of an invertible matrix A , we required that A be square and that the inverse matrix B satisfy *both* $AB = I_n$ and $BA = I_n$. In fact, given that A and B are square,

either of the conditions $AB = I_n$ and $BA = I_n$ implies the other one. For example, assume $A, B \in M_n(F)$ and $AB = I_n$. Taking determinants gives $\det(A)\det(B) = \det(I_n) = 1_F$. Since F is commutative, we also have $1_F = \det(B)\det(A)$, so that $\det(A)$ is an invertible element of F . By one of the theorems quoted in the last paragraph, A^{-1} (the unique two-sided inverse of A) exists. Multiplying both sides of $AB = I_n$ on the left by A^{-1} , we get $B = A^{-1}$, so that $BA = A^{-1}A = I_n$. Similarly, one proves that $BA = I_n$ implies $AB = I_n$. However, our proof (which used determinants) relies very heavily on the fact that the matrices in question are square. For $A \in M_{m,n}(F)$, we say $B \in M_{n,m}(F)$ is a *left inverse* of A iff $BA = I_n$; we say $B \in M_{n,m}(F)$ is a *right inverse* of A iff $AB = I_m$. For $m \neq n$, one can give examples where B is a left inverse but not a right inverse of A , and vice versa.

We can now prove that for $U, V \in M_n(F)$, U and V are invertible iff UV is invertible. We proved the forward implication already by checking that $(UV)^{-1} = V^{-1}U^{-1}$. Conversely, assume UV is invertible with inverse $W \in M_n(F)$. Then $U(VW) = (UV)W = I_n$, so the theorem proved in the last paragraph shows that VW is the inverse of U . Similarly, $(WU)V = W(UV) = I_n$ shows that WU is the inverse of V . More generally, induction on k shows that for $U_1, \dots, U_k \in M_n(F)$, every U_i is invertible iff the product $U_1U_2 \cdots U_k$ is invertible.

4.7 Matrix Operations via Columns

As we have seen, any $m \times n$ matrix A can be identified with the list of its columns in $(F^m)^n$: $A = (A^{[1]}, A^{[2]}, \dots, A^{[n]})$, where $A^{[j]} = (A(1, j), A(2, j), \dots, A(m, j))$ is column j of A . In symbols, $A^{[j]}(i) = A(i, j)$ for $1 \leq i \leq m$ and $1 \leq j \leq n$. It is fruitful to recast the definitions of the matrix operations in terms of the columns of the matrices involved.

First, assume $A, B \in M_{m,n}(F)$ and $c \in F$. We have $(A + B)^{[j]} = A^{[j]} + B^{[j]} \in F^m$ since, for $1 \leq i \leq m$, $(A + B)^{[j]}(i) = (A + B)(i, j) = A(i, j) + B(i, j) = A^{[j]}(i) + B^{[j]}(i) = (A^{[j]} + B^{[j]})(i)$. One sees similarly that $(cA)^{[j]} = c \cdot A^{[j]} \in F^m$. This means that we can add two matrices of the same size column by column, and multiplying a matrix by a scalar multiplies each column by that scalar.

Next, assume $A \in M_{m,n}(F)$ and $B \in M_{n,p}(F)$, and let $C = AB \in M_{m,p}(F)$. How are the columns of A and B related to the columns of C ? To answer this, let $1 \leq i \leq m$ and $1 \leq j \leq p$, and compute

$$C^{[j]}(i) = (AB)(i, j) = \sum_{k=1}^n A(i, k)B(k, j) = \sum_{k=1}^n A(i, k)B^{[j]}(k) = (A(B^{[j]}))(i).$$

(The last equality is the defining formula for multiplying the matrix A on the right by the column vector $B^{[j]}$.) So $(AB)^{[j]} = A(B^{[j]})$ for all j , which means that *the j 'th column of AB can be found by multiplying the matrix A by the j 'th column of B* . For example, the second column of the product

$$\begin{bmatrix} 2 & 0 \\ -1 & 3 \\ 0 & 4 \end{bmatrix} \begin{bmatrix} \pi & 2 & e & 3i & \sqrt{37} \\ 9! & 3 & \sqrt[3]{-11} & 10^9 & (2+3i)^7 \end{bmatrix}$$

is

$$\begin{bmatrix} 2 & 0 \\ -1 & 3 \\ 0 & 4 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 4 \\ 7 \\ 12 \end{bmatrix}.$$

Now consider the product Av , where A is an $m \times n$ matrix and v is an $n \times 1$ column vector (or element of F^n). For all $i \in [m]$, we have

$$(Av)(i) = \sum_{j=1}^n A(i, j)v(j) = \sum_{j=1}^n A^{[j]}(i)v(j) = \left(\sum_{j=1}^n A^{[j]}v(j) \right) (i).$$

So $Av = \sum_{j=1}^n A^{[j]}v(j)$ in F^m [or $Av = \sum_{j=1}^n v(j)A^{[j]}$ for commutative rings F], which means that the matrix-vector product Av is a linear combination of the columns of A — namely, the linear combination whose coefficients are the entries in v . For example,

$$\begin{bmatrix} 1 & 2 & 0 & -1 \\ 3 & i & -2 & 4 \\ 1/2 & 0 & 0 & \pi \end{bmatrix} \begin{bmatrix} 2 \\ -1 \\ 5 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ -i - 4 \\ 1 \end{bmatrix} = 2 \begin{bmatrix} 1 \\ 3 \\ 1/2 \end{bmatrix} - 1 \begin{bmatrix} 2 \\ i \\ 0 \end{bmatrix} + 5 \begin{bmatrix} 0 \\ -2 \\ 0 \end{bmatrix} + 0 \begin{bmatrix} -1 \\ 4 \\ \pi \end{bmatrix}.$$

[The rest of this section applies only to fields F .] Let us deduce some consequences of the fact that Av is a linear combination of the columns of A . Given $A \in M_{m,n}(F)$, define the range of A to be the set $R(A) = \{Av : v \in F^n\} \subseteq F^m$, which is readily seen to be a subspace of F^m . Our formula for Av shows that

$$R(A) = \{v(1)A^{[1]} + \cdots + v(n)A^{[n]} : v(i) \in F\},$$

so that the range of A coincides with the subspace of all linear combinations of the columns of A in F^m . This subspace is also called the *column space* of A and may be denoted $\text{Col}(A)$. The dimension of this subspace is called the *column rank* of A , denoted $\text{colrk}(A)$.

For another application, consider again the general matrix product AB , where $A \in M_{m,n}(F)$ and $B \in M_{n,p}(F)$. We will show that $\text{colrk}(AB) \leq \min(\text{colrk}(A), \text{colrk}(B))$, or (equivalently) $\text{colrk}(AB) \leq \text{colrk}(A)$ and $\text{colrk}(AB) \leq \text{colrk}(B)$. First, since $(AB)^{[j]} = A(B^{[j]})$, where $B^{[j]}$ is a column vector, we see that the j 'th column of AB is a linear combination of the columns of A , where the coefficients come from the j 'th column of B . In particular, each column $(AB)^{[j]}$ belongs to the column space of A . Knowing that $(AB)^{[j]} \in \text{Col}(A)$ for all j , it follows that $\text{Col}(AB) \subseteq \text{Col}(A)$ since the columns of AB generate the column space of AB . Taking dimensions, we see that $\text{colrk}(AB) \leq \text{colrk}(A)$. Now let us show that $\text{colrk}(AB) \leq \text{colrk}(B)$. Define a map $L_A : F^n \rightarrow F^m$ by sending each column vector $v \in F^n$ to the column vector $L_A(v) = Av \in F^m$. One checks that L_A is a linear map. Let $W = \text{Col}(B) \subseteq F^n$, and recall that W is spanned by the columns $B^{[1]}, \dots, B^{[p]}$ of B . Applying the linear map L_A , we see (using Exercise 45 of Chapter 1) that $L_A[W]$ is spanned by the vectors $L_A(B^{[j]}) = A(B^{[j]}) = (AB)^{[j]}$, for $j = 1, \dots, p$. The latter vectors are exactly the columns of AB . We conclude that $\text{Col}(AB) = L_A[W]$. Since the column space of AB is the image of the column space of B under a linear map, we must have $\dim(\text{Col}(AB)) \leq \dim(\text{Col}(B))$ (see Exercise 43 in Chapter 1). In other words, $\text{colrk}(AB) \leq \text{colrk}(B)$, as needed.

Now suppose $A \in M_{m,n}(F)$. If P is an invertible $m \times m$ matrix, we have inequalities

$$\text{colrk}(A) \geq \text{colrk}(PA) \geq \text{colrk}(P^{-1}(PA)) = \text{colrk}(A),$$

so that $\text{colrk}(PA) = \text{colrk}(A)$. Similarly, if Q is an invertible $n \times n$ matrix, the inequalities

$$\text{colrk}(A) \geq \text{colrk}(AQ) \geq \text{colrk}((AQ)Q^{-1}) = \text{colrk}(A)$$

show that $\text{colrk}(AQ) = \text{colrk}(A)$. Thus, multiplication on the left or right by an invertible matrix does not change the column rank. In particular, letting A be the identity matrix I_n , which has column rank n since its columns constitute the basis (e_1, \dots, e_n) of F^n , we see that every invertible $n \times n$ matrix has column rank n .

4.8 Matrix Operations via Rows

We now recast the definitions of matrix operations in terms of the rows of the matrices and vectors involved. This time, we start with the fact that any $m \times n$ matrix A can be identified with the list of its rows: $A = (A_{[1]}, A_{[2]}, \dots, A_{[m]}) \in (F^n)^m$, where $A_{[i]} = (A(i, 1), A(i, 2), \dots, A(i, n))$ is row i of A . In symbols, $A_{[i]}(j) = A(i, j)$ for $1 \leq i \leq m$ and $1 \leq j \leq n$.

As in the case of columns, a short calculation shows that $(A + B)_{[i]} = A_{[i]} + B_{[i]}$ and $(cA)_{[i]} = c(A_{[i]})$ in F^n for all $A, B \in M_{m,n}(F)$, all $c \in F$, and all $i \in [m]$. This means that we can add two matrices of the same size row by row, and multiplying a matrix by a scalar multiplies each row by that scalar.

Next, assume $A \in M_{m,n}(F)$ and $B \in M_{n,p}(F)$, and let $C = AB \in M_{m,p}(F)$. How are the rows of A and B related to the rows of C ? Answer: for any $1 \leq i \leq m$ and $1 \leq j \leq p$, we have

$$C_{[i]}(j) = (AB)(i, j) = \sum_{k=1}^n A(i, k)B(k, j) = \sum_{k=1}^n A_{[i]}(k)B(k, j) = (A_{[i]}B)(j).$$

So $(AB)_{[i]} = (A_{[i]})B$ for all $i \in [m]$, which means that *the i 'th row of AB can be found by multiplying the i 'th row of A by the matrix B* . For example, the third row of the product

$$\begin{bmatrix} 2 & 0 \\ -1 & 3 \\ 0 & 4 \end{bmatrix} \begin{bmatrix} 5 & 7 & 9 & 11 & 13 \\ 1 & 1 & 1 & 1 & 10 \end{bmatrix}$$

is

$$\begin{bmatrix} 0 & 4 \end{bmatrix} \begin{bmatrix} 5 & 7 & 9 & 11 & 13 \\ 1 & 1 & 1 & 1 & 10 \end{bmatrix} = \begin{bmatrix} 4 & 4 & 4 & 4 & 40 \end{bmatrix}.$$

Now consider the product wA , where A is an $m \times n$ matrix and w is a $1 \times m$ row vector (or element of F^m). For all $j \in [n]$, we have

$$(wA)(j) = \sum_{i=1}^m w(i)A(i, j) = \sum_{i=1}^m w(i)A_{[i]}(j) = \left(\sum_{i=1}^m w(i)A_{[i]} \right) (j).$$

So $wA = \sum_{i=1}^m w(i)A_{[i]}$ in F^n [and the scalars $w(i)$ must appear on the left for non-commutative rings F]. In words, *the vector-matrix product wA is a linear combination of the rows of A — namely, the linear combination whose coefficients are the entries in the row vector w* . For example,

$$\begin{bmatrix} 3 & 0 & -2 \end{bmatrix} \begin{bmatrix} 1 & 2 & 0 \\ 3 & i & -2 \\ 1/2 & 0 & 1 \end{bmatrix} = 3 \begin{bmatrix} 1 & 2 & 0 \end{bmatrix} + 0 \begin{bmatrix} 3 & i & -2 \end{bmatrix} - 2 \begin{bmatrix} 1/2 & 0 & 1 \end{bmatrix}.$$

[The rest of this section applies only to *fields* F .] For $A \in M_{m,n}(F)$, define the *row space* of A to be the subspace $\text{Row}(A) \subseteq F^n$ spanned by the m rows of A , which consists of all linear combinations of the rows of A . Define the *row rank* of A , denoted $\text{rowrk}(A)$, to be the dimension of the row space of A . Arguing as we did earlier for columns, our formula for wA shows that $\text{Row}(A) = \{wA : w \in F^m\} \subseteq F^n$. Furthermore, for $B \in M_{n,p}(F)$, each row i of the matrix product AB has the form $(AB)_{[i]} = (A_{[i]})B$, hence is a linear combination of the rows of B with coefficients coming from row i of A . Since the rows of AB generate $\text{Row}(AB)$,

we conclude that $\text{Row}(AB) \subseteq \text{Row}(B)$ and $\text{rowrk}(AB) \leq \text{rowrk}(B)$. On the other hand, consider the linear map $R_B : F^n \rightarrow F^p$ that sends a row vector $w \in F^n$ to the row vector $R_B(w) = wB \in F^p$. Let $V = R_B[\text{Row}(A)]$ be the image of $\text{Row}(A)$ under this linear map. Since $\text{Row}(A)$ is spanned by $A_{[1]}, \dots, A_{[m]}$, V is spanned by $R_B(A_{[1]}), \dots, R_B(A_{[m]})$. In other words, V is spanned by $A_{[1]}B = (AB)_{[1]}, \dots, A_{[m]}B = (AB)_{[m]}$. But these are exactly the rows of AB , and hence $V = \text{Row}(AB)$. Since $V = R_B[\text{Row}(A)]$, we must have $\dim(V) \leq \dim(\text{Row}(A))$, and hence $\text{rowrk}(AB) \leq \text{rowrk}(A)$. In summary,

$$\text{rowrk}(AB) \leq \min(\text{rowrk}(A), \text{rowrk}(B)).$$

Once again, fix $A \in M_{m,n}(F)$. If P is an invertible $m \times m$ matrix, we have inequalities

$$\text{rowrk}(A) \geq \text{rowrk}(PA) \geq \text{rowrk}(P^{-1}(PA)) = \text{rowrk}(A),$$

so that $\text{rowrk}(PA) = \text{rowrk}(A)$. Similarly, if Q is an invertible $n \times n$ matrix, the inequalities

$$\text{rowrk}(A) \geq \text{rowrk}(AQ) \geq \text{rowrk}((AQ)Q^{-1}) = \text{rowrk}(A)$$

show that $\text{rowrk}(AQ) = \text{rowrk}(A)$. Thus, *multiplication on the left or right by an invertible matrix does not change the row rank*. In particular, letting A be the identity matrix I_n , which has row rank n since its rows comprise the basis (e_1, \dots, e_n) of F^n , we see that *every invertible $n \times n$ matrix has row rank n* . Later in this chapter, we will prove that $\text{rowrk}(A) = \text{colrk}(A)$ for every matrix A .

Before leaving this section, let us mention interpretations for the transpose and conjugate-transpose operations in terms of rows and columns. Suppose A is an $m \times n$ matrix. From the definition $A^T(i, j) = A(j, i)$, we see immediately that $(A^T)_{[i]}(j) = A^T(i, j) = A(j, i) = A^{[i]}(j)$ for $1 \leq j \leq m$, so that $(A^T)_{[i]} = A^{[i]}$ for all $i \in [n]$. Similarly, $(A^T)_{[j]}(i) = A^T(i, j) = A(j, i) = A_{[j]}(i)$ for all $i \in [n]$ implies that $(A^T)_{[j]} = A_{[j]}$ for all $j \in [m]$. Translating these equations into words, we see that: the i 'th row of A^T is the i 'th column of A (for all $i \in [n]$), while the j 'th column of A^T is the j 'th row of A (for all $j \in [m]$). Thus, the transposition operation “interchanges the rows and columns of A .” An analogous calculation when $F = \mathbb{C}$ shows that $(A^*)_{[i]} = \overline{A^{[i]}}$ and $(A^*)_{[j]} = \overline{A_{[j]}}$, where the complex conjugate of a row or column vector is found by conjugating each entry. So, we obtain A^* from A by interchanging rows and columns and taking the complex conjugate of every entry.

4.9 Elementary Operations and Elementary Matrices

This section discusses the link between elementary row or column operations on matrices and multiplication by certain matrices called *elementary matrices*. Recall from introductory linear algebra that there are three types of *elementary row operations* that we can apply to an $m \times n$ matrix A . First, for any $i \in [m]$, we can multiply the i 'th row of A by some *nonzero* scalar c in the field F . [In the case of rings, we require c to be an invertible element (unit) in the ring F .] For non-commutative rings, each entry of the row is multiplied *on the left* by c .] Second, for all $i \neq j$ in $[m]$, we can interchange the i 'th row and the j 'th row of A . Third, for all $i \neq j$ in $[m]$ and any scalar $c \in F$, we can add c times the j 'th row to the i 'th row of A [c is on the left for non-commutative rings]. There are three analogous *elementary column operations* that act on the columns of A . [For non-commutative rings, c multiplies column entries *on the right*.] These elementary operations are helpful for solving

systems of linear equations, for computing normal forms and other invariants of matrices, and for other applications.

We now define three kinds of *elementary matrices*. First, for $i \in [p]$ and nonzero $c \in F$ [or units c in the case of rings], let $E_p^1[c; i]$ be the $p \times p$ matrix with i, i -entry equal to c , j, j -entries equal to 1_F for $j \neq i$, and all other entries zero. Second, for $i \neq j$ in $[p]$, let $E_p^2[i, j]$ be the $p \times p$ matrix with ones in positions (i, j) , (j, i) , and (k, k) for all $k \neq i, j$, and zeroes elsewhere. Third, for $c \in F$ and $i \neq j$ in $[p]$, let $E_p^3[c; i, j]$ be the $p \times p$ matrix with ones on the diagonal, a c in position (i, j) , and zeroes elsewhere. As we will see shortly, the three types of elementary matrices encode the three types of elementary row or column operations defined above. As a mnemonic aid, note that one can obtain the elementary matrix for a given elementary operation by performing that operation on the $p \times p$ identity matrix. For example,

$$E_3^1[4; 2] = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 1 \end{bmatrix}; \quad E_3^2[2, 3] = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}; \quad E_3^3[3; 2, 1] = \begin{bmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

Now fix an $m \times n$ matrix A , and take $p = m$ above. We claim that *multiplying A on the left by one of the elementary matrices has the same effect as applying the corresponding elementary row operation to A* . This is not hard to check by writing down formulas for all the entries, but it is more instructive to give a conceptual proof using the ideas from §4.8. For example, consider the matrix $B = E_m^3[c; i, j]A$. The k 'th row of B is $E_m^3[c; i, j]_{[k]}A$. In turn, this is the linear combination of the m rows of A with coefficients given by the entries in the row vector $E_m^3[c; i, j]_{[k]}$. If $k \neq i$, the row vector in question is e_k , so the row $B_{[k]}$ is 1 times $A_{[k]}$ plus 0 times every other row of A ; i.e., the k 'th row of B equals the k 'th row of A . If $k = i$, $E_m^3[c; i, j]_{[i]} = ce_j + e_i$ by definition, so the row $B_{[i]}$ is c times $A_{[j]}$ plus 1 times $A_{[i]}$; i.e., the i 'th row of B is c times the j 'th row of A plus the i 'th row of A . This proves that multiplication on the left by $E_m^3[c; i, j]$ has the same effect as the elementary row operation of the third type. The assertions for $E_m^1[c; i]$ and $E_m^2[i, j]$ are checked in the same way.

On the other hand, *if we multiply $A \in M_{m,n}(F)$ on the right by an $n \times n$ elementary matrix, this has the same effect as applying the corresponding elementary column operation to A* . For example, consider the matrix $C = AE_n^2[i, j]$. For $k \in [n]$, the k 'th column of C is $A(E_n^2[i, j]^{[k]})$, which is the linear combination of the columns of A with coefficients given by the entries in the column vector $E_n^2[i, j]^{[k]}$. If $k \neq i, j$, this column vector is e_k , so $C^{[k]} = A^{[k]}$. If $k = i$, this column vector is e_j , so $C^{[i]} = A^{[j]}$. If $k = j$, this column vector is e_i , so $C^{[j]} = A^{[i]}$. So the columns of C are precisely the columns of A with column i and column j interchanged. The assertions for the other elementary matrices are checked in the same way. Note carefully that to add c times row j to row i , one left-multiplies by $E_m^3[c; i, j]$; but, to add c times column j to column i , one right-multiplies by $E_n^3[c; j, i]$.

One may check that *the inverse of an elementary matrix exists and is an elementary matrix of the same type*. Specifically, $E_p^1[c; i]^{-1} = E_p^1[c^{-1}; i]$, $E_p^2[i, j]^{-1} = E_p^2[i, j]$, and $E_p^3[c; i, j]^{-1} = E_p^3[-c; i, j]$. One can verify these formulas by computing with entries, but they can also be checked by applying the remarks in the previous two paragraphs. For example, the product $E_p^3[c; i, j]E_p^3[-c; i, j]$ is the matrix obtained from $E_p^3[-c; i, j]$ by adding c times row j to row i . Row j of $E_p^3[-c; i, j]$ is e_j , and row i is $e_i + (-c)e_j$, so the new row i after the row operation is e_i . For all $k \neq i$, row k of $E_p^3[-c; i, j]$ is e_k , and this is unchanged by the row operation. So $E_p^3[c; i, j]E_p^3[-c; i, j] = I_p$, and the product in the other order is the identity for analogous reasons. The formulas for the inverses of the other types of elementary matrices can be similarly verified.

4.10 Elementary Matrices and Gaussian Elimination

[In the rest of the chapter, we assume F is a field.] The next several sections derive some results whose proofs use elementary matrices. Another fundamental fact that can be proved with elementary matrices is the product formula $\det(AB) = \det(A)\det(B)$ for $A, B \in M_n(F)$, which we will prove in §5.9.

Let us begin by reviewing some basic facts concerning the *Gaussian elimination* algorithm (which we assume the reader has seen before) and relating this algorithm to elementary matrices. For a more detailed study of Gaussian elimination and its connection to matrix factorizations, see Chapter 9. Starting with any $m \times n$ matrix A , the Gaussian elimination algorithm applies a sequence of elementary row operations to transform A to its *reduced row-echelon form*, denoted A_{ech} . To describe A_{ech} conveniently, let us call the leftmost nonzero entry in a given nonzero row of A_{ech} the *leading entry* of that row. Then A_{ech} is characterized by the following properties: all zero rows of A_{ech} (if any) occur at the bottom of the matrix; the leading entry of each nonzero row is 1_F ; for all $i \in [m - 1]$, the leading entry in row i is strictly left of the leading entry in row $i + 1$ (if any); and for each leading entry, the other elements in the column containing that entry are all zero. It can be proved that for every $A \in M_{m,n}(F)$, there exists a unique reduced row-echelon form A_{ech} that can be reached from A by applying finitely many elementary row operations. Existence is not hard to prove using Gaussian elimination (more precisely, Gauss–Jordan elimination, which produces the zeroes in the columns of the leading entries). The uniqueness of A_{ech} is more subtle, but we will not need to invoke uniqueness in the following discussion.

We have seen that each elementary row operation used in the passage from A to A_{ech} can be accomplished by multiplying on the left by an appropriate elementary matrix. Thus, $A_{\text{ech}} = E_k E_{k-1} \cdots E_1 A$, where E_s is the elementary matrix corresponding to the s 'th row operation applied by the Gaussian elimination algorithm.

For example, let us compute the reduced row-echelon form for the matrix

$$A = \begin{bmatrix} 0 & 3 & -4 & 1 \\ 2 & 4 & -2 & 6 \\ 4 & -1 & 8 & 9 \end{bmatrix}$$

and find a specific factorization $A_{\text{ech}} = E_k \cdots E_1 A$. To obtain a leading 1 in the 1,1-position, we switch row 1 and row 2 of A and then multiply the new row 1 by $1/2$. This can be accomplished by left-multiplying A first by $E_1 = E_3^2[1, 2]$, and then by $E_2 = E_3^1[1/2; 1]$:

$$E_1 A = \begin{bmatrix} 2 & 4 & -2 & 6 \\ 0 & 3 & -4 & 1 \\ 4 & -1 & 8 & 9 \end{bmatrix}; \quad E_2 E_1 A = \begin{bmatrix} 1 & 2 & -1 & 3 \\ 0 & 3 & -4 & 1 \\ 4 & -1 & 8 & 9 \end{bmatrix}.$$

Next we add -4 times row 1 to row 3, by left-multiplying by the elementary matrix $E_3 = E_3^3[-4; 3, 1]$; and we continue by adding 3 times row 2 to row 3 (left-multiply by $E_4 = E_3^3[3; 3, 2]$):

$$E_3 E_2 E_1 A = \begin{bmatrix} 1 & 2 & -1 & 3 \\ 0 & 3 & -4 & 1 \\ 0 & -9 & 12 & -3 \end{bmatrix}; \quad E_4 E_3 E_2 E_1 A = \begin{bmatrix} 1 & 2 & -1 & 3 \\ 0 & 3 & -4 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

To finish, we multiply row 2 by $1/3$ (left-multiply by $E_5 = E_3^1[1/3; 2]$) and then add -2

times the new row 2 to row 1 (left-multiply by $E_6 = E_3^3[-2; 1, 2]$):

$$E_5 E_4 E_3 E_2 E_1 A = \begin{bmatrix} 1 & 2 & -1 & 3 \\ 0 & 1 & -4/3 & 1/3 \\ 0 & 0 & 0 & 0 \end{bmatrix};$$

$$A_{\text{ech}} = E_6 E_5 E_4 E_3 E_2 E_1 A = \begin{bmatrix} 1 & 0 & 5/3 & 7/3 \\ 0 & 1 & -4/3 & 1/3 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

4.11 Elementary Matrices and Invertibility

Our first theorem involving elementary matrices states that *a square matrix A is invertible iff A can be written as a product of elementary matrices*. To prove this, assume first that A is a product of elementary matrices. We know every elementary matrix is invertible, so any product of elementary matrices is also invertible (§4.6). Conversely, assume that $A \in M_n(F)$ is invertible. Write $A_{\text{ech}} = E_k E_{k-1} \cdots E_1 A$ for certain elementary matrices E_j . Since A and each elementary matrix E_j is invertible, so is their product A_{ech} . It readily follows that A_{ech} cannot have any zero rows. Since A_{ech} is square, the definition of reduced row-echelon form forces A_{ech} to be the identity matrix I_n . Consequently, $E_k E_{k-1} \cdots E_1 A = I_n$. Solving for A , we find that $A = E_1^{-1} E_2^{-1} \cdots E_k^{-1}$. Since the inverse of an elementary matrix is elementary, A has been expressed as a product of elementary matrices.

For example, let us express the invertible matrix $A = \begin{bmatrix} 2 & 1 \\ 4 & 3 \end{bmatrix}$ as a product of elementary matrices. We reduce A to $A_{\text{ech}} = I_2$ by the following sequence of elementary row operations: add -2 times row 1 to row 2, producing $\begin{bmatrix} 2 & 1 \\ 0 & 1 \end{bmatrix}$; multiply row 1 by $1/2$, producing $\begin{bmatrix} 1 & 1/2 \\ 0 & 1 \end{bmatrix}$; add $-1/2$ times row 2 to row 1, producing I_2 . In terms of matrices, we have

$$I_2 = E_2^3[-1/2; 1, 2]E_2^1[1/2; 1]E_2^3[-2; 2, 1]A = \begin{bmatrix} 1 & -1/2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1/2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ 4 & 3 \end{bmatrix}.$$

Solving for A gives

$$\begin{bmatrix} 2 & 1 \\ 4 & 3 \end{bmatrix} = A = E_2^3[2; 2, 1]E_2^1[2; 1]E_2^3[1/2; 1, 2] = \begin{bmatrix} 1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1/2 \\ 0 & 1 \end{bmatrix}.$$

4.12 Row Rank and Column Rank

As our second application of elementary matrices, we prove that *for any $m \times n$ matrix A , $\text{rowrk}(A) = \text{colrk}(A)$* . We can therefore define the *rank of A* , denoted $\text{rank}(A)$, as the common value of the row rank of A and the column rank of A . Recall that we have already shown (§4.8) that multiplying a matrix on the left or right by an invertible matrix does not change the row rank; we also proved that these operations do not affect the column rank either (§4.7). Moreover, we have seen (§4.10) that $A_{\text{ech}} = PA$ for some matrix P that must be invertible (being a product of elementary matrices).

Let us continue to simplify the matrix A_{ech} by applying elementary *column* operations, which can be effected by multiplying on the right by certain elementary matrices. The leading 1 in each nonzero row is the only nonzero entry *in its column*. Adding appropriately chosen constant multiples of this column to the columns to its right, we can arrange that the leading 1 in each nonzero row is the only nonzero entry *in its row*. Our matrix now consists of some number r of nonzero rows, each of which has a single 1 in it, and these 1's appear in distinct columns. Using column interchanges, we can arrange the matrix so that these 1's appear on the first r positions of the main diagonal. (The *main diagonal* of a rectangular $m \times n$ matrix consists of the i, i -entries for $1 \leq i \leq \min(m, n)$.) The final matrix is obtained from $A_{\text{ech}} = PA$ by right-multiplying by some matrix Q , which is a product of elementary matrices and is therefore invertible.

In summary, given $A \in M_{m,n}(F)$, we can find invertible matrices $P \in M_m(F)$ and $Q \in M_n(F)$ and an integer $r \in \{0, 1, \dots, \min(m, n)\}$ such that $(PAQ)_{i,i} = 1$ for $1 \leq i \leq r$, and all other entries of PAQ are zero. The matrix P (resp. Q) is the product of the elementary matrices corresponding to the elementary row operations (resp. column operations) we applied to transform A into the indicated form. The matrix PAQ is called the *Smith canonical form* or *Smith normal form* for A . [See Chapter 18 for a discussion of the more general case in which F is a principal ideal domain (PID).] Now we can quickly prove our main result. For, we know that $\text{rowrk}(A) = \text{rowrk}(PAQ)$ and $\text{colrk}(A) = \text{colrk}(PAQ)$. But it is evident from the form of PAQ that $\text{rowrk}(PAQ) = r = \text{colrk}(PAQ)$. Thus, $\text{rowrk}(A) = \text{colrk}(A)$. We remark that the proof provides an algorithm for computing the rank of A . One can also show that the row rank (and hence the column rank) of A equals the number of nonzero rows in the reduced row-echelon form A_{ech} . This observation gives a faster way to compute the rank of A .

For example, the matrix $A = \begin{bmatrix} 0 & 3 & -4 & 1 \\ 2 & 4 & -2 & 6 \\ 4 & -1 & 8 & 9 \end{bmatrix}$ from §4.10 has reduced row-echelon form

$$A_{\text{ech}} = PA = \begin{bmatrix} 1 & 0 & 5/3 & 7/3 \\ 0 & 1 & -4/3 & 1/3 \\ 0 & 0 & 0 & 0 \end{bmatrix},$$

where $P = \begin{bmatrix} -2/3 & 1/2 & 0 \\ 1/3 & 0 & 0 \\ 3 & -2 & 1 \end{bmatrix}$ is the product of six elementary matrices used to reduce A to A_{ech} .

It is already evident that $\text{rowrk}(A) = \text{rowrk}(A_{\text{ech}}) = 2$ (and it is almost as quick to determine that $\text{colrk}(A) = \text{colrk}(A_{\text{ech}}) = 2$). However, if we wanted to eliminate the nonzero entries to the right of the leading 1's in A_{ech} , we could do so by the following elementary column operations: add $-5/3$ times column 1 to column 3; add $-7/3$ times column 1 to column 4; add $4/3$ times column 2 to column 3; add $-1/3$ times column 2 to column 4. In terms of matrices, if we define

$$Q = E_4^3[-5/3; 1, 3]E_4^3[-7/3; 1, 4]E_4^3[4/3; 2, 3]E_4^3[-1/3; 2, 4] = \begin{bmatrix} 1 & 0 & -5/3 & -7/3 \\ 0 & 1 & 4/3 & -1/3 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix},$$

then $PAQ = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$ has the form prescribed in the proof above, with $r = 2$ nonzero entries on the main diagonal.

4.13 Conditions for Invertibility of a Matrix

Table 4.1 displays one of the central theorems of matrix algebra, which gives a long list of conditions that are all logically equivalent to the invertibility of a square matrix A . Condition (8b) says that the system of linear equations $Ax = b$ always has a solution x for any right-hand side b ; condition (10b) says the system always has a unique solution for any b ; and condition (9b) says the homogeneous system $Ax = 0$ has no solutions other than the zero solution $x = 0$. As we will see, these conditions can be reformulated to say that a certain linear map L_A (left multiplication by A) is surjective, or bijective, or injective. To detect when all these conditions hold, it is sufficient to check that A has nonzero determinant, or to test the rows (or columns) of A for linear independence, or to see if the reduced row-echelon form of A is I_n . We also see from conditions (2a) and (2b) that A will be invertible provided it has either a left inverse or a right inverse.

TABLE 4.1

Conditions Equivalent to the Invertibility of A .

Assume F is a field, $A \in M_n(F)$, and $L_A : F^n \rightarrow F^n$ is the map given by $L_A(x) = Ax$ for all column vectors $x \in F^n$. The following conditions are all equivalent:

- | | |
|-------|---|
| (1) | A is invertible. |
| (2a) | $BA = I_n$ for some $B \in M_n(F)$. |
| (2b) | $AC = I_n$ for some $C \in M_n(F)$. |
| (3) | $\det(A) \neq 0_F$. |
| (4) | A is a product of finitely many elementary matrices. |
| (5) | The reduced row-echelon form A_{ech} is I_n . |
| (6a) | $\text{rowrk}(A) = n$. |
| (6b) | The n rows of A form a linearly independent list in F^n . |
| (7a) | $\text{colrk}(A) = n$. |
| (7b) | The n columns of A form a linearly independent list in F^n . |
| (7c) | The range of A has dimension n . |
| (8a) | The map L_A is surjective (onto). |
| (8b) | For all $b \in F^n$, there exists $x \in F^n$ with $Ax = b$. |
| (9a) | The map L_A is injective (one-to-one). |
| (9b) | For all $x \in F^n$, if $Ax = 0$ then $x = 0$. |
| (9c) | The null space of A is $\{0\}$. |
| (10a) | The map L_A is bijective (and hence an isomorphism). |
| (10b) | For all $b \in F^n$, there exists a unique $x \in F^n$ with $Ax = b$. |
| (11) | A^T is invertible. |
| (12) | [when $F = \mathbb{C}$] A^* is invertible. |

We now turn to the proof of the theorem in Table 4.1, which involves the verification of many conditional and biconditional statements.

(1) \Rightarrow (2a): If A is invertible, then (2a) certainly holds by choosing $B = A^{-1}$.

(1) \Rightarrow (2b): If A is invertible, then (2b) certainly holds by choosing $C = A^{-1}$.

(2a) \Rightarrow (3): Given $BA = I_n$, take determinants and use the product formula to get $\det(B)\det(A) = \det(I_n) = 1_F$. Since $1_F \neq 0_F$ in a field F , this forces $\det(A) \neq 0_F$.

(2b) \Rightarrow (3): The proof is analogous to (2a) \Rightarrow (3). (Similarly, (1) implies (3).)

(3) \Rightarrow (1): If $\det(A) \neq 0_F$, then we can write down an explicit formula for A^{-1} (see §5.11) to confirm that A has a two-sided inverse. This formula involves a division by $\det(A)$, which is why we need A to have nonzero determinant.

(1) \Leftrightarrow (4): We proved this in §4.11.

(1) \Leftrightarrow (5): For any matrix A , we saw in §4.10 that $A_{\text{ech}} = PA$, where $P \in M_n(F)$ is a product of elementary matrices, hence is invertible. So, A is invertible iff A_{ech} is invertible (§4.6). Let $B \in M_n(F)$ be any matrix in reduced row-echelon form. Using the definition of reduced row-echelon form, we see that B has no zero rows iff every row of B contains a leading 1 iff every column of B contains a leading 1 iff $B = I_n$. Since I_n is invertible and no matrix with a row of zeroes is invertible, we see that B is invertible iff $B = I_n$. Applying this remark to $B = A_{\text{ech}}$, we see that A is invertible iff $A_{\text{ech}} = I_n$.

(1) \Leftrightarrow (6a): As above, $A_{\text{ech}} = PA$ for some invertible matrix P . Multiplying by an invertible matrix does not change the row rank, so $\text{rowrk}(A_{\text{ech}}) = \text{rowrk}(A)$. If A is invertible, then $A_{\text{ech}} = I_n$ (since (1) implies (5)), so $\text{rowrk}(A_{\text{ech}}) = n$, so (6a) is true. Conversely, if (6a) holds, then $\text{rowrk}(A_{\text{ech}}) = n$, so A_{ech} cannot have any zero rows. As shown in the last paragraph, this forces $A_{\text{ech}} = I_n$. Since (5) implies (1), A is invertible.

(6a) \Leftrightarrow (6b): By definition, the row space of A is spanned by the n rows of A . If (6b) holds, then the list of n rows of A is a basis for the row space of A , so that $\text{rowrk}(A) = n$. Conversely, if (6b) fails, then we can delete one or more vectors from the list of rows to get a basis for the row space, forcing $\text{rowrk}(A) < n$.

(6a) \Leftrightarrow (7a): This follows from the identity $\text{rowrk}(A) = \text{colrk}(A)$, valid for all matrices A . We proved this identity in §4.12.

(7a) \Leftrightarrow (7b) \Leftrightarrow (7c): We prove the equivalence of (7a) and (7b) in the same way we proved the equivalence of (6a) and (6b). The range of A is the same as the column space of A (§4.7), so (7a) is equivalent to (7c) by definition.

(7c) \Leftrightarrow (8a): The range of the matrix A is the subspace $\{Ax : x \in F^n\} = \{L_A(x) : x \in F^n\}$, which is precisely the image of the map L_A . Note $L_A : F^n \rightarrow F^n$ is surjective iff the image of L_A is all of F^n iff the dimension of the image is n (since $\dim(F^n) = n$, but all proper subspaces of F^n have dimension smaller than n) iff the range of A has dimension n .

(8a) \Leftrightarrow (8b): Once we recall that $Ax = L_A(x)$, we see that (8b) is the very definition of L_A being surjective.

(8a) \Leftrightarrow (9c): Let the linear map $L_A : F^n \rightarrow F^n$ have image R and kernel N . We know R is the range of A . Moreover, N is the null space of A , since $x \in N$ iff $L_A(x) = 0$ iff $Ax = 0$. By the rank-nullity theorem (see §1.8), $\dim(R) + \dim(N) = \dim(F^n) = n$. Now (8a) is equivalent to (7c), which says $\dim(R) = n$. In turn, this is equivalent to $\dim(N) = 0$, which is equivalent to $N = \{0\}$.

(9a) \Leftrightarrow (9c): The linear map L_A is injective iff its kernel is $\{0\}$. As observed above, the kernel of L_A is the same as the null space of A .

(9b) \Leftrightarrow (9c): The null space of A is $\{0\}$ iff the null space of A is a subset of $\{0\}$. Writing out the definition of the latter condition, we obtain (9b).

(9a) \Leftrightarrow (10a): Assume (9a) holds. Then (8a) holds also, so L_A is injective and surjective, so (10a) holds. The converse is immediate, since all bijections are one-to-one.

(10a) \Leftrightarrow (10b): Recalling that $Ax = L_A(x)$, we see that (10b) is none other than the definition of the bijectivity of L_A .

(1) \Leftrightarrow (11),(1) \Leftrightarrow (12): We proved these equivalences in §4.6.

4.14 Summary

1. *Formal Definitions for Matrices.* An element of F^n is a function $x : [n] \rightarrow F$, which is often presented as a list $(x(1), \dots, x(n))$. An $m \times n$ matrix is a function $A : [m] \times [n] \rightarrow F$; special cases are column vectors ($n = 1$) and row vectors

($m = 1$), which can be identified with elements in F^m and F^n , respectively. A matrix can be identified with a list of its columns ($A^{[1]}, \dots, A^{[n]}$) or with a list of its rows ($A_{[1]}, \dots, A_{[m]}$); by definition, $A^{[j]}(i) = A(i, j) = A_{[i]}(j)$ for $i \in [m]$ and $j \in [n]$.

2. *Vector Spaces of Functions.* If S is any set and F is a field, the set V of all functions from S to F becomes an F -vector space under pointwise addition and scalar multiplication of functions. For $x \in S$, let $e_x \in V$ be the function that sends x to 1_F and everything else in S to zero. When S is finite, the set $X = \{e_x : x \in S\}$ is a basis for V such that $|X| = |S|$. In particular, $\dim(V) = |S|$ in this case.
3. *Spaces of Vectors and Matrices.* The sets F^n and $M_{m,n}(F)$ are F -vector spaces under componentwise operations. Viewing vectors and matrices as functions, these operations are the same as pointwise addition and scalar multiplication of functions from $[n]$ or $[m] \times [n]$ into F . A basis for F^n consists of the vectors e_i defined by $e_i(i) = 1$ and $e_i(j) = 0$ for $j \neq i$. A basis for $M_{m,n}(F)$ consists of the matrices e_{ij} defined by $e_{ij}(i, j) = 1$ and $e_{ij}(i', j') = 0$ whenever $i \neq i'$ or $j \neq j'$. We have $\dim(F^n) = n$ and $\dim(M_{m,n}(F)) = mn$.
4. *Matrix Multiplication.* For $A \in M_{m,n}(F)$ and $B \in M_{n,p}(F)$, $AB \in M_{m,p}(F)$ is the matrix whose i, j -entry is $\sum_{k=1}^n A(i, k)B(k, j)$. We have $AI_n = A = I_m A$ and $c(AB) = (cA)B = A(cB)$ for $c \in F$ [and F commutative]. Matrix multiplication is distributive and associative. When $m = n$ [and F is commutative], the product operation turns the vector space $M_n(F)$ into an F -algebra. Generalized associativity holds for products of more than three terms, which can therefore be written with no parentheses. If A_u has order $n_{u-1} \times n_u$ for all u , the i, j -entry of $\prod_{u=1}^s A_u$ is
$$\sum_{k_1=1}^{n_1} \sum_{k_2=1}^{n_2} \cdots \sum_{k_{s-1}=1}^{n_{s-1}} A_1(i, k_1)A_2(k_1, k_2) \cdots A_{s-1}(k_{s-2}, k_{s-1})A_s(k_{s-1}, j).$$
5. *Transpose and Conjugate-Transpose.* For $A \in M_{m,n}(F)$, we define $A^T \in M_{n,m}(F)$ and [when $F = \mathbb{C}$] $A^* \in M_{n,m}(\mathbb{C})$ by $A^T(i, j) = A(j, i)$ and $A^*(i, j) = \overline{A(j, i)}$ for $i \in [n]$ and $j \in [m]$. These operations satisfy the following identities: $(A + B)^T = A^T + B^T$; $(cA)^T = c(A^T)$ for $c \in F$; $(A_1 \cdots A_s)^T = A_s^T \cdots A_1^T$ [for F commutative]; $(A + B)^* = A^* + B^*$; $(cA)^* = \bar{c}(A^*)$ for $c \in \mathbb{C}$; and $(A_1 \cdots A_s)^* = A_s^* \cdots A_1^*$.
6. *Significance of Rows and Columns.* On one hand, $(AB)^{[j]} = A(B^{[j]})$, which says that the j 'th column of AB is A times the j 'th column of B . On the other hand, $(AB)_{[i]} = (A_{[i]}B)$, which says that the i 'th row of AB is the i 'th row of A times B . For a matrix-vector product Av , we can write $Av = \sum_j A^{[j]}v(j)$, so that Av is a linear combination of the columns of A . For a vector-matrix product wA , we can write $wA = \sum_i w(i)A_{[i]}$, so that wA is a linear combination of the rows of A . Some consequences [for fields F] are: $\text{colrk}(AB) \leq \min(\text{colrk}(A), \text{colrk}(B))$; $\text{rowrk}(AB) \leq \min(\text{rowrk}(A), \text{rowrk}(B))$; left or right multiplication by an invertible matrix does not change the row rank or the column rank; and the rank of an invertible $n \times n$ matrix is equal to n . Matrix transposition interchanges the rows and columns of A ; the conjugate-transpose operation also conjugates every entry.
7. *Elementary Operations vs. Elementary Matrices.* The three elementary row operations are: interchange two rows; multiply a row by an invertible scalar; add any scalar multiple of one row to a different row. To apply an elementary

row operation to A , one need only multiply A on the left by the corresponding elementary matrix (formed by applying the same row operation to the identity matrix). There are analogous elementary column operations, which can be achieved by multiplying A on the right by the corresponding elementary matrix. Elementary matrices are invertible; the inverses are also elementary.

8. *Applications of Elementary Matrices.* [This item applies when F is a field.] A square matrix is invertible iff it is a product of elementary matrices. For any A (possibly rectangular), there exists an invertible matrix P such that PA is in reduced row-echelon form. Furthermore, there exists another invertible matrix Q such that PAQ is a matrix in Smith canonical form, with $r = \text{rank}(A)$ ones on the main diagonal and all other entries zero. In particular, the column rank of A is always equal to the row rank of A . Elementary matrices can be used to prove the product formula $\det(AB) = \det(A)\det(B)$ for determinants.
9. *Invertible Matrices.* [This item applies when F is a field.] A matrix $A \in M_n(F)$ is invertible iff there exists a (necessarily unique) matrix $A^{-1} \in M_n(F)$ with $AA^{-1} = I_n = A^{-1}A$. The following conditions are equivalent to the existence of A^{-1} : A has a left inverse; A has a right inverse; A has nonzero determinant; A is a product of elementary matrices; A 's reduced row-echelon form is I_n ; A has row rank n ; A has column rank n ; the n rows of A are linearly independent; the n columns of A are linearly independent; the range of A has dimension n ; the null space of A is $\{0\}$; the map $L_A : F^n \rightarrow F^n$ given by $L_A(x) = Ax$ for $x \in F^n$ is injective; L_A is surjective; L_A is bijective; the system $Ax = b$ has a solution $x \in F^n$ for every $b \in F^n$; $Ax = b$ has a *unique* solution x for every b ; $Ax = 0$ has only the zero solution $x = 0$; A^T is invertible; [when $F = \mathbb{C}$] A^* is invertible. If A^{-1} exists, then $(A^{-1})^{-1} = A$, $(A^T)^{-1} = (A^{-1})^T$, and [when $F = \mathbb{C}$] $(A^*)^{-1} = (A^{-1})^*$. A product $A = A_1 A_2 \cdots A_s$ is invertible iff every A_i is invertible, in which case $A^{-1} = A_s^{-1} \cdots A_2^{-1} A_1^{-1}$. $\text{GL}_n(F)$, the set of invertible matrices in $M_n(F)$, is a group under matrix multiplication.

4.15 Exercises

Unless otherwise stated, assume F is a field in these exercises.

1. (a) State a precise formula for a map $I : F^n \rightarrow M_{n,1}(F)$ that identifies n -tuples with column vectors, and check that I is a vector space isomorphism. (Note that the inputs and outputs for I are functions.) (b) Repeat (a) for a map $J : F^n \rightarrow M_{1,n}(F)$ that identifies n -tuples with row vectors.
2. (a) State a precise formula for a map $I : F \rightarrow F^1$, and prove that I is a vector space isomorphism. (b) State a precise formula for a map $J : F \rightarrow M_{1,1}(F)$, and prove that J is a vector space isomorphism. (This problem allows us to identify 1-tuples and 1×1 matrices with scalars in F .)
3. (a) In the text, we noted there is a map $C : M_{m,n}(F) \rightarrow (F^m)^n$ sending $A \in M_{m,n}(F)$ to its list of columns $(A^{[1]}, \dots, A^{[n]})$. Give a precise formula for C^{-1} , and check that C is F -linear. (b) Repeat (a) for the map $R : M_{m,n}(F) \rightarrow (F^n)^m$ sending $A \in M_{m,n}(F)$ to its list of rows $(A_{[1]}, \dots, A_{[m]})$.
4. Let S be a set, and let V be the set of all functions $g : S \rightarrow F$ with pointwise addition and scalar multiplication. (a) Check that addition on V is commutative.

- (b) Check that the zero function from S to F is the additive identity element of V . (c) Check that the additive inverse of $f \in V$ is the function sending each $x \in S$ to $-f(x)$. (d) Check that scalar multiplication on V satisfies the five axioms in the definition of a vector space.
5. Let S be any set, and let E be any F -vector space. Let V be the set of all functions $g : S \rightarrow E$. (a) Define “pointwise” addition and scalar multiplication operations on V , and prove that V with these operations is an F -vector space. (b) Suppose $S = \{x_1, \dots, x_k\}$ is finite, and E is finite-dimensional with ordered basis (v_1, \dots, v_m) . Describe an explicit basis for V , and compute $\dim_F(V)$.
6. (a) Cite a theorem from the text to explain why the set V of all functions $f : \mathbb{R} \rightarrow \mathbb{R}$ (under pointwise operations) is a real vector space. (b) For each $x \in \mathbb{R}$, let $e_x : \mathbb{R} \rightarrow \mathbb{R}$ send x to 1 and all other inputs to 0. Which functions in V are in the span of the set $\{e_x : x \in \mathbb{R}\}$? Is this set a basis for V ? (c) Explain why the set of continuous $f : \mathbb{R} \rightarrow \mathbb{R}$ is a subspace of V . Find the intersection of this subspace and the span of $\{e_x : x \in \mathbb{R}\}$.
7. Let $A = \begin{bmatrix} 5 & 2 & 1 \\ 1/2 & -3 & 0 \\ -1 & -4 & -3/2 \end{bmatrix}$, $B = \begin{bmatrix} 1 & -1 & 0 \\ 0 & 2 & -2 \\ 3 & -3 & 0 \end{bmatrix}$, $C = \begin{bmatrix} 4 & 2 & -7 \\ 1 & 1/3 & 3 \end{bmatrix}$, $w = [2 \ 4 \ 1]$, and $v = [1 \ -3 \ -1]$. Compute: (a) $A + B$; (b) AB ; (c) BA ; (d) CA ; (e) CB ; (f) $A(v^T)$; (g) wB ; (h) $w(v^T)$ (i) v^Tw .
8. Let $A = \begin{bmatrix} 5-i & 2i+3 \\ 0 & -1-3i \end{bmatrix}$, $B = \begin{bmatrix} 3+4i & -2i \\ e^{5\pi i/3} & \sqrt{2}+2i \end{bmatrix}$, and $v = \begin{bmatrix} i \\ 1 \end{bmatrix}$. Compute: (a) $A + B$; (b) AB ; (c) A^T ; (d) B^T ; (e) A^* ; (f) B^* ; (g) Bv ; (h) v^Tv ; (i) v^*v ; (j) v^*Av ; (k) v^*A^*Av .
9. (a) Prove: for $A, B \in M_n(F)$, A and B commute iff A^T and B^T commute. (b) Prove: for invertible $A, B \in M_n(F)$, A and B commute iff A^{-1} and B^{-1} commute.
10. Let $A = \begin{bmatrix} 1 & 4 \\ 2 & 3 \end{bmatrix}$. (a) Assuming $A \in M_2(\mathbb{R})$, find A^2 , A^3 , AA^T , and (if possible) A^{-1} . (b) Repeat (a), assuming $A \in M_2(\mathbb{Z}_5)$. (c) Repeat (a), assuming $A \in M_2(\mathbb{Z}_7)$.
11. Let $F = \mathbb{Z}_2$, $A = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$, $B = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$, and $C = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix}$. (a) Compute A^k for all $k \in \mathbb{Z}$. (b) Compute B^k for all $k \in \mathbb{Z}$. (c) Compute C^k for all $k \in \mathbb{Z}$.
12. Let $A, A' \in M_{m,n}(F)$, $B \in M_{n,p}(F)$, and $c \in F$. Prove carefully: (a) $(A+A')B = AB + A'B$; (b) $c(AB) = (cA)B$; (c) $c(AB) = A(cB)$; (d) $AI_n = A$. (e) Let F be a non-commutative ring. Must all of (a) through (d) be true? Explain.
13. Let $A, B \in M_{m,n}(\mathbb{C})$ and $C \in M_{n,p}(\mathbb{C})$. Prove carefully: (a) $(A+B)^* = A^* + B^*$; (b) $(AC)^* = C^*A^*$; (c) $(A^*)^* = A$; (d) A is invertible iff A^* is invertible, in which case $(A^*)^{-1} = (A^{-1})^*$.
14. *Hadamard Product of Matrices.* For matrices $A, B \in M_{m,n}(F)$, define their *Hadamard product* to be the matrix $A \odot B \in M_{m,n}(F)$ given by $(A \odot B)(i, j) = A(i, j)B(i, j)$ for $i \in [m]$ and $j \in [n]$. We obtain $A \odot B$ by multiplying corresponding entries of A and B . Prove that \odot : (a) is commutative; (b) is associative; (c) satisfies the left and right distributive laws with respect to

- matrix addition; (d) has an identity (what is it?). (e) Which matrices in $M_{m,n}(F)$ have inverses with respect to \odot ?
15. (a) Give a specific example of matrices $A, B \in M_n(F)$ with $AB = 0$ and $BA \neq 0$.
 (b) Give a specific example of matrices $A, B \in M_{m,n}(F)$ with $m \neq n$, $BA = 0$, and $AB \neq 0$. (c) If possible, find matrices $A, B \in M_4(\mathbb{R})$ such that every entry of A and B is nonzero, yet $AB = BA = 0$.
 16. A matrix $A \in M_n(F)$ is called *diagonal* iff for all $i, j \in [n]$, $i \neq j$ implies $A(i, j) = 0_F$. This means that the only nonzero entries of A occur on the main diagonal. Prove that the set of diagonal matrices is a commutative subalgebra of the F -algebra $M_n(F)$. Find a basis for this subalgebra and compute its dimension.
 17. A matrix $A \in M_n(F)$ is called *upper-triangular* iff for all $i, j \in [n]$, $i > j$ implies $A(i, j) = 0_F$. This means that all entries of A below the main diagonal are zero. Let $A, B \in M_n(F)$ be upper-triangular matrices. (a) Prove $A + B$ is upper-triangular. (b) Prove cA is upper-triangular for all $c \in F$. (c) Is 0_n upper-triangular? Is I_n upper-triangular? Is J_n upper-triangular? (d) Prove AB is upper-triangular.
 18. A matrix $A \in M_n(F)$ is called *lower-triangular* iff for all $i, j \in [n]$, $i < j$ implies $A(i, j) = 0_F$. This means that all entries of A above the main diagonal are zero. Prove A is lower-triangular iff A^T is upper-triangular. Then prove analogues of (a) through (d) in Exercise 17 for lower-triangular matrices. Give very short proofs by appealing to properties of matrix transpose.
 19. A matrix $A \in M_n(F)$ is called *strictly upper-triangular* iff for all $i, j \in [n]$, $i \geq j$ implies $A(i, j) = 0_F$. This means that the only nonzero entries of A occur strictly above the main diagonal. Assume A is strictly upper-triangular.
 (a) Prove $A^n = 0_n$ using the formula for $[\prod_{u=1}^s A_u] (i, j)$ in §4.5. Deduce that A is not invertible. (b) Prove: if $B \in M_n(F)$ is upper-triangular, then AB and BA are strictly upper-triangular. (c) Show that the set of strictly upper-triangular matrices is a subspace of $M_n(F)$. Exhibit a basis for this subspace and compute its dimension.
 20. A matrix $A \in M_n(F)$ is called *unitriangular* iff A is upper-triangular and $A(i, i) = 1_F$ for all $i \in [n]$. Let S be the set of unitriangular matrices in $M_n(F)$. Determine if S is closed under: (a) addition; (b) additive inverses; (c) additive identity; (d) multiplication; (e) multiplicative inverses; (f) multiplicative identity; (g) scalar multiplication. [Exercise 42 below may help with (e).]
 21. A matrix $A \in M_n(F)$ is called *nilpotent* iff $A^k = 0$ for some integer $k > 0$. For example, you showed in Exercise 19(a) that any strictly upper-triangular matrix is nilpotent. (a) Give an example of a nilpotent matrix that is neither upper-triangular nor lower-triangular. (b) Give an example of two nilpotent matrices whose sum is not nilpotent. (c) Give an example of two nilpotent matrices whose product is not nilpotent. (d) Prove: if $A, B \in M_n(F)$ are nilpotent and $AB = BA$, then $A + B$ and AB are nilpotent.
 22. A matrix $A \in M_n(F)$ is called *skew-symmetric* iff $A^T = -A$. (a) Prove that the set of skew-symmetric matrices is a subspace of $M_n(F)$, and compute its dimension. (b) Give an example of a nonzero matrix that is both symmetric and skew-symmetric, or explain why no such matrix exists. (c) Prove: if $A, B \in M_n(F)$ are skew-symmetric and commute, then AB is symmetric. Does the result hold without assuming $AB = BA$?
 23. The *trace* of a square matrix $A \in M_n(F)$ is defined by $\text{tr}(A) = \sum_{i=1}^n A(i, i)$,

which is the sum of the main diagonal entries of A . Let $A, B \in M_n(F)$. (a) Compute $\text{tr}(0_n)$, $\text{tr}(I_n)$, and $\text{tr}(J_n)$. (b) Prove: $\text{tr}(A+B) = \text{tr}(A) + \text{tr}(B)$ and $\text{tr}(cA) = c\text{tr}(A)$ for all $c \in F$. (c) Prove: $\text{tr}(AB) = \text{tr}(BA)$. Does this identity hold if $A \in M_{m,n}(F)$ and $B \in M_{n,m}(F)$ with $n \neq m$? (d) Prove: $\text{tr}(A^T) = \text{tr}(A)$. For $F = \mathbb{C}$, find an identity relating $\text{tr}(A^*)$ and $\text{tr}(A)$. (e) Prove or disprove: for invertible A , $\text{tr}(A^{-1}) = \text{tr}(A)^{-1}$. (f) Give an example of a field F and a matrix $A \in M_n(F)$ with all entries of A nonzero and $\text{tr}(A) = 0_F$.

24. *Commutator of Matrices.* For all $A, B \in M_n(F)$, define the *commutator* of A and B to be $[A, B] = AB - BA$. (a) Prove: for all $A, B \in M_n(F)$, $[A, B] = 0$ iff A and B commute. (b) Prove: for all $A, B \in M_n(F)$, $[A, B] = -[B, A]$ (anticommutativity). (c) Prove: for all $A, B, C \in M_n(F)$ and all $d \in F$, $[A+B, C] = [A, C] + [B, C]$, $[A, B+C] = [A, B] + [A, C]$, and $[dA, B] = d[A, B] = [A, dB]$ (bilinearity). (d) Prove: for all $A, B, C \in M_n(F)$, $[A, [B, C]] + [B, [C, A]] + [C, [A, B]] = 0$ (Jacobi identity). (e) Prove: for all $A, B \in M_n(F)$, the trace of $[A, B]$ is zero. (f) Is $[A, [B, C]] = [[A, B], C]$ true for all $A, B, C \in M_n(F)$? Either prove it, or find conditions on A, B, C that will cause this identity to hold.
25. Write out all possible complete parenthesizations of a product $ABCDE$ of five matrices.
26. For $n \geq 1$, let p_n be the number of complete parenthesizations of a product of n matrices; so, for example, $p_1 = p_2 = 1$, $p_3 = 2$, and $p_4 = 5$. Show that $p_n = \sum_{k=1}^{n-1} p_k p_{n-k}$. (It can be shown from this recursion that p_n is the *Catalan number* $\frac{1}{n}\binom{2n-2}{n-1}$.)
27. *Laws of Exponents for Matrix Powers.* Let $A \in M_n(F)$. (a) Prove: for all integers $r, s \geq 0$, $A^{r+s} = A^r A^s$. (b) Prove: for all integers $r, s \geq 0$, $(A^r)^s = A^{rs}$. (c) Show that (a) and (b) hold for possibly negative integers r and s , provided that A^{-1} exists.
28. Let (G, \star) be any group, $s \geq 1$, and $a_1, \dots, a_s \in G$. State and prove a generalized associativity result that justifies writing the product $a_1 \star a_2 \star \dots \star a_s$ with no parentheses. (Imitate the proof given in the text for matrix products.)
29. (a) Suppose $U \in M_{m,n}(F)$ and $V \in M_{n,p}(F)$. How many additions and multiplications in F are required to compute UV using the definition of matrix product? (b) Suppose A is 2×10 , B is 10×15 , C is 15×4 , and D is 4×8 . For each complete parenthesization of $ABCD$, count the number of multiplications needed to evaluate this matrix product by multiplying matrices in the order determined by the parenthesization.
30. (a) Prove: for all $a, b, c, d \in F$, if $ad - bc \neq 0_F$, then $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ is invertible, and $A^{-1} = (ad - bc)^{-1} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$. (b) Assuming $ad - bc \neq 0_F$, find an explicit factorization of A into a product of elementary matrices.
31. Find the inverse of $A = \begin{bmatrix} 0 & 1 & 2 & 1 \\ 1 & 0 & 1 & 1 \\ 0 & 2 & 2 & 1 \\ 1 & 2 & 2 & 1 \end{bmatrix}$ in $M_4(\mathbb{Z}_3)$.
32. Find necessary and sufficient conditions on the integer n and the field F for the group $\text{GL}_n(F)$ to be commutative.
33. (a) Find the size of the group $\text{GL}_3(\mathbb{Z}_2)$. [Hint: Build $A \in \text{GL}_3(\mathbb{Z}_2)$ by choosing nonzero rows, one at a time, that are not in the span of the previous rows.]

- (b) Find the size of the group $\mathrm{GL}_2(\mathbb{Z}_7)$. (c) For any $n \geq 1$ and any prime p , find a formula for the size of $\mathrm{GL}_n(\mathbb{Z}_p)$.
34. (a) Give a specific example of functions $f, g : \mathbb{R} \rightarrow \mathbb{R}$ with $f \circ g = \mathrm{id}_{\mathbb{R}}$ but $g \circ f \neq \mathrm{id}_{\mathbb{R}}$. (b) Is f onto? Is f one-to-one? Is g onto? Is g one-to-one? (c) Does there exist another example in (a) for which the answers to (b) would be different?
35. (a) Show that $A = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \end{bmatrix} \in M_{2,3}(\mathbb{R})$ has a right inverse but no left inverse.
Find all right inverses of A . (b) Give a specific example of a matrix $A \in M_{4,3}(\mathbb{R})$ that has a left inverse but no right inverse. (c) Is the left inverse of the matrix in (b) unique? If not, can you find an example where the left inverse is unique?
36. (a) Prove: for $A \in M_{m,n}(F)$, $\mathrm{rank}(A) \leq \min(m, n)$. (b) Prove: for $A, B \in M_{m,n}(F)$, $\mathrm{rank}(A + B) \leq \mathrm{rank}(A) + \mathrm{rank}(B)$. Give an example where strict inequality holds.
37. Let $m < n$ be positive integers. By studying rank, show that no $A \in M_{m,n}(F)$ can have a left inverse, and no $B \in M_{n,m}(F)$ can have a right inverse.
38. Given $A \in M_{m,n}(F)$, $B \in M_{n,p}(F)$, $i \in [m]$, and $j \in [p]$, prove that $[(AB)(i,j)] = A_{[i]} B^{[j]}$. Give a verbal description of what this equation means.
39. Define $A \in M_{3,5}(\mathbb{R})$ and $B \in M_{5,3}(\mathbb{R})$ by setting $A(i,j) = i^2 j$ and $B(j,i) = j + 2i$ for $i \in [3]$ and $j \in [5]$. (a) Compute row 3 of AB . (b) Compute column 2 of AB . (c) Express column 4 of BA as a linear combination of the columns of B . (d) Express row 1 of BA as a linear combination of the rows of A .
40. Let $A, B \in M_{m,n}(F)$ and $c \in F$. Prove, as stated in the text, that $(A + B)_{[i]} = A_{[i]} + B_{[i]}$ and $(cA)_{[i]} = c(A_{[i]})$ for all $i \in [m]$.
41. (a) Let $A \in M_n(F)$ be upper-triangular with $A(i,i) \neq 0$ for all $i \in [n]$, and let $B \in M_{m,n}(F)$ be arbitrary. Prove: for all $i \in [n]$, the subspace of F^n spanned by the first i columns of B equals the subspace of F^n spanned by the first i columns of BA . (b) Formulate and prove a similar result involving left-multiplication by an upper-triangular matrix with no zeroes on the diagonal. (c) Formulate and prove results similar to (a) and (b) involving multiplication on the left or right by lower-triangular matrices with no zeroes on the diagonal.
42. (a) Let $A \in M_n(F)$ be unitriangular (see Exercise 20). Prove A is invertible, and find a recursive formula for computing the entries of A^{-1} . [Hint: Solve $BA = I_n$ for the columns of B from left to right.] (b) Illustrate your formula by finding the inverses of

$$A = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 2 & 3 \\ 0 & 0 & 1 & 1 & 2 \\ 0 & 0 & 0 & 1 & 3 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \text{ and } C = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}.$$

43. Which lower-triangular $A \in M_n(F)$ are invertible? Prove that A^{-1} , when it exists, is also lower-triangular.
44. Given $A \in M_n(F)$ and positive integers $r < s$, prove with minimal calculation that $\mathrm{Col}(A^s) \subseteq \mathrm{Col}(A^r)$; similarly, prove $\mathrm{Row}(A^s) \subseteq \mathrm{Row}(A^r)$.

45. For each matrix below, determine the rank of the matrix by finding a subset of the columns that forms a basis for the column space: (a) the matrix B in Exercise 7; (b) the matrix C in Exercise 7; (c) $J_{m,n}$, which has every entry equal to 1;

$$(d) \text{ the matrix } A = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} \in M_4(\mathbb{Z}_2); (e) \text{ the matrix } A \in M_{m,n}(\mathbb{R})$$

with $A(i,j) = ij$ for $i \in [m]$ and $j \in [n]$; (f) the matrix $A \in M_n(\mathbb{R})$ with $A(i,j) = (i-1)n + j$ for $i, j \in [n]$.

46. Repeat Exercise 45, but compute the rank by finding a subset of the rows that forms a basis for the row space.
47. Repeat Exercise 45, but compute the rank by finding the reduced row-echelon form of the matrix and counting the nonzero rows.
48. Let $A \in M_{m,n}(F)$, $i, \ell \in [m]$, $j, k \in [n]$, and $c, d \in F$ with $c \neq 0$. (a) Prove that left-multiplication by $E_m^1[c; i]$ multiplies row i of A by c . (b) Prove that right-multiplication by $E_n^1[c; j]$ multiplies column j of A by c . (c) Prove $E_m^1[c; i]$ is invertible with inverse $E_m^1[c^{-1}; i]$. (d) Prove that for $i \neq \ell$, left-multiplication by $E_m^2[i, \ell]$ interchanges row i and row ℓ of A . (e) Prove that for $i \neq \ell$, $E_m^2[i, \ell]$ is its own inverse. (f) Prove that for $j \neq k$, right-multiplication by $E_n^3[d; j, k]$ adds d times column j of A to column k of A .
49. Prove that every elementary matrix $E_p^2[i, j]$ can be expressed as a product of several elementary matrices of the form $E_p^3[\pm 1; r, s]$. Conclude that the elementary operation of interchanging two rows of a matrix can be simulated by other elementary row operations.
50. (a) List all elementary matrices in $M_3(\mathbb{Z}_2)$. (b) Express each matrix in (a) as a product $D_1 D_2 \cdots D_k$, where each D_i is one of the matrices A, B, C displayed in Exercise 11.
51. Draw pictures of all possible matrices $A \in M_{3,5}(F)$ that are in reduced row-echelon form. Each matrix entry should be a zero, a one, or a star to indicate an arbitrary scalar.
52. For each matrix A below, compute the reduced row-echelon form A_{ech} and an invertible matrix P (a product of elementary matrices) with $A_{\text{ech}} = PA$. (a) $A \in M_{3,4}(\mathbb{R})$ given by $A(i,j) = i$ for $i \in [3]$ and $j \in [4]$; (b) $A \in M_4(\mathbb{R})$ given by $A(i,j) = (-1)^{i+j}$ for $i, j \in [4]$; (c) A in Exercise 8;

$$(d) A = \begin{bmatrix} 2 & 4 & -3 & 0 & -14 \\ 1 & 2 & 4 & -2 & 13 \\ 0 & 0 & 1 & 0 & 4 \\ 2 & 4 & 2 & 0 & 6 \end{bmatrix}; (e) A = \begin{bmatrix} 0 & 1 & 0 & 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 \end{bmatrix} \in M_{8,7}(\mathbb{Z}_2).$$

53. For each matrix A in Exercise 52, find an invertible matrix Q (a product of elementary matrices) such that PAQ is the Smith canonical form of A .
54. If possible, write each matrix below and its inverse as a product of elementary matrices. (a) A in Exercise 7; (b) B in Exercise 7; (c) A in Exercise 31; (d) A in Exercise 42(b).

55. Prove: for $A, B \in M_{m,n}(F)$, $\text{rank}(A) = \text{rank}(B)$ iff there exist invertible matrices $U \in M_m(F)$ and $V \in M_n(F)$ with $A = UBV$.
56. In Table 4.1, use the definitions to give a direct proof of (7b) \Leftrightarrow (9b).
57. In Table 4.1, assume (1), (6b), and (11) are known to be equivalent. Using this, prove (6b) \Leftrightarrow (7b).
58. In Table 4.1, prove (8b) implies (2b) by solving for the columns of C one at a time.
59. Let $A \in M_n(F)$, and define $R_A : F^n \rightarrow F^n$ by $R_A(x) = xA$ for each row vector $x \in F^n$. Prove the following conditions are equivalent: (a) A is invertible; (b) R_A is injective; (c) R_A is surjective; (d) R_A is bijective.
60. Let $T : V \rightarrow W$ be an F -linear map between two n -dimensional F -vector spaces V and W , where $n \in \mathbb{N}$. Prove the following conditions are equivalent: (a) T is injective; (b) T is surjective; (c) T is bijective.
61. (a) Give an example of an F -vector space V and a linear map $T : V \rightarrow V$ that is one-to-one but not onto. (b) Give an example of an F -vector space V and a linear map $S : V \rightarrow V$ that is onto but not one-to-one. (c) Give an example of an F -vector space V and linear maps $T, S : V \rightarrow V$ with $S \circ T = \text{id}_V$ but $T \circ S \neq \text{id}_V$.
62. Fix positive integers $m < n$ and a matrix $A \in M_{m,n}(F)$. Prove that the following conditions on A are all equivalent: (a) A has a right inverse. (b) A has at least two right inverses. (c) The reduced row-echelon form of A has no zero rows. (d) $\text{rowrk}(A) = m$. (e) The list of the m rows of A is linearly independent in F^n . (f) $\text{colrk}(A) = m$. (g) The range of A has dimension m . (h) The map $L_A : F^n \rightarrow F^m$ sending each $x \in F^n$ to Ax is surjective. (i) For all $b \in F^m$, there exists $x \in F^n$ with $Ax = b$. (j) A^T has a left inverse.
63. Fix positive integers $m > n$ and a matrix $A \in M_{m,n}(F)$. Formulate and prove a list of conditions on A (similar to those in Exercise 62) that are all equivalent to A having a left inverse.
64. This exercise studies proofs of the theorem: *for all $A, B \in M_n(F)$, if $BA = I_n$ then $AB = I_n$* . (a) Show that the theorem follows from the statement: for all $A, B \in M_n(F)$, if $BA = I_n$ then A is invertible. (b) We proved (a) in the text using determinants (see the proofs of (2a) implies (3) implies (1) in §4.13). Give several new proofs of (a) by using the equivalences proved in §4.13 (ignoring conditions (2a) and (2b)) and showing that (2a) implies (6a), (2a) implies (7a), and (2a) implies (9b). (c) Similarly, give proofs that $AC = I_n$ implies A is invertible by showing (2b) implies (6a), (2b) implies (7a), (2b) implies (8b), and by using (11) and (a) of this exercise. (d) Prove, without using determinants, that for all $U, V \in M_n(F)$, if UV is invertible then U and V are invertible.
65. This exercise outlines another proof [37] that $\text{rowrk}(A) = \text{colrk}(A)$ for all $A \in M_{m,n}(F)$. Let us call a row of A *extraneous* iff the row is an F -linear combination of the other rows of A ; define *extraneous columns* analogously. (a) Prove: if A has no extraneous rows, then $\text{rowrk}(A) = m$ and $m \leq n$. (b) Prove: if A has no extraneous columns, then $\text{colrk}(A) = n$ and $n \leq m$. (c) Explain why $\text{rowrk}(A) = \text{colrk}(A)$ if A has no extraneous rows and no extraneous columns. (d) Suppose A has an extraneous row. Argue that deletion of this row produces a smaller matrix with the same row rank *and* column rank as A . [*Hint:* If row i is extraneous, consider the map $T : F^m \rightarrow F^{m-1}$ that deletes the i 'th coordinate. Show that the restriction of T to the column space of A is injective.] (e) Prove a result similar to (d) if A has an extraneous column. (f) Show that repeated use of

- (d) and (e) always leads to a smaller matrix B having no extraneous rows or columns, such that $\text{rowrk}(B) = \text{rowrk}(A)$ and $\text{colrk}(B) = \text{colrk}(A)$. Explain why $\text{rowrk}(A) = \text{colrk}(A)$ follows.
66. This exercise outlines yet another proof, due to Hans Liebeck [35], that $\text{colrk}(A) = \text{rowrk}(A)$ for all $A \in M_n(\mathbb{C})$. (a) Prove: for all $y \in M_{n,1}(\mathbb{C})$, $y = 0$ iff $y^*y = 0$. (b) Prove: for all $A \in M_n(\mathbb{C})$ and $x \in M_{n,1}(\mathbb{C})$, $Ax = 0$ iff $A^*Ax = 0$. (c) Prove: for all $A \in M_n(\mathbb{C})$, $\text{colrk}(A) = \text{colrk}(A^*A)$. (d) Prove: for all $A \in M_n(\mathbb{C})$, $\text{colrk}(A^*) = \text{rowrk}(A)$. (e) Prove: for all $A \in M_n(\mathbb{C})$, $\text{colrk}(A) \leq \text{colrk}(A^*)$. (f) Explain how to deduce $\text{colrk}(A) = \text{rowrk}(A)$ from (d) and (e).

This page intentionally left blank

Determinants via Calculations

In introductory linear algebra courses, students are told the basic properties of determinants and are given formulas for computing and manipulating determinants. However, such courses rarely provide complete justifications for these properties and formulas. For example, the determinant of a general $n \times n$ matrix A is often defined recursively by Laplace expansion along the first row. We are told that using Laplace expansion along *any* row or column gives the same answer, but this claim is seldom proved.

In this chapter, we rigorously define determinants by an explicit, non-recursive formula that involves a sum indexed by permutations. (The background on permutations needed here is covered in Chapter 2.) We use the explicit formula to prove some familiar properties of determinants, including the aforementioned Laplace expansions. We give computational proofs whenever possible. These proofs work for matrices with entries in any commutative ring (as defined in §1.2), so we operate at that level of generality. Chapter 20 introduces a more sophisticated approach to determinants based on multilinear algebra. While that approach is conceptually harder, it yields many properties of determinants with a minimum of calculation.

5.1 Matrices with Entries in a Ring

Before defining determinants, we review the definitions of matrices with entries in a ring and operations on such matrices. (This material is developed more fully in Chapter 4.) Let R be a ring. For each positive integer m , let $[m] = \{1, 2, \dots, m\}$. An $m \times n$ *matrix* A over R is a function $A : [m] \times [n] \rightarrow R$. For $i \in [m]$ and $j \in [n]$, the ring element $A(i, j)$ is frequently denoted $A_{i,j}$ or A_{ij} or a_{ij} , and we often display the matrix A as an array of m rows and n columns such that $A(i, j)$ appears in row i and column j . Let $M_{m,n}(R)$ be the set of all $m \times n$ matrices with entries in R , and let $M_n(R)$ be the set of all $n \times n$ (square) matrices with entries in R . For each $n \geq 1$, we have a *zero matrix* $0_n \in M_n(R)$ defined by setting $0_n(i, j) = 0_R$ for all $i, j \in [n]$. For each $n \geq 1$, we also define the *identity matrix* I_n by setting $I_n(i, i) = 1_R$ for all $i \in [n]$ and $I_n(i, j) = 0_R$ for all $i \neq j$ in $[n]$. We write 0 for 0_n and I for I_n if n is understood from context.

A *row vector* is an element of $M_{1,n}(R)$; a *column vector* is an element of $M_{n,1}(R)$. We often identify row vectors and column vectors with n -tuples in R^n , which are defined to be functions $x : [n] \rightarrow R$. We often represent such a function by the ordered list $(x(1), \dots, x(n))$. If x is a row vector, we may write $x(i)$ or x_i instead of $x(1, i)$; similarly if x is a column vector. We prefer using parentheses instead of subscripts here, since we will often be dealing with lists of row vectors. For instance, if (x_1, \dots, x_m) is a list of row vectors, then $x_3(5)$ is the fifth entry in the row vector x_3 .

Table 5.1 defines some algebraic operations on matrices. The basic properties of these operations were covered in Chapter 4. Note that the operations \bar{A} and A^* are defined only for matrices with complex entries.

TABLE 5.1

Definitions of Matrix Operations. Here R is a ring, $A, B \in M_{m,n}(R)$, $C \in M_{n,p}(R)$, $d \in R$, and bars denote complex conjugation ($\overline{x+iy} = x-iy$ for $x, y \in \mathbb{R}$).

Matrix Operation	Defining Formula
1. matrix sum $A+B \in M_{m,n}(R)$	$(A+B)(i,j) = A(i,j) + B(i,j)$
2. scalar multiplication $dA \in M_{m,n}(R)$	$(dA)(i,j) = d \cdot A(i,j)$
3. matrix product $BC \in M_{m,p}(R)$	$(BC)(i,j) = \sum_{k=1}^n B(i,k) \cdot C(k,j)$
4. transpose $A^T \in M_{n,m}(R)$	$A^T(i,j) = A(j,i)$
5. matrix conjugate $\bar{A} \in M_{m,n}(\mathbb{C})$	$\bar{A}(i,j) = \overline{A(i,j)}$
6. conjugate-transpose $A^* \in M_{n,m}(\mathbb{C})$	$A^*(i,j) = \overline{A(j,i)}$

5.2 Explicit Definition of the Determinant

From now on, R will be a fixed *commutative* ring. The *determinant* is a function that associates to each $n \times n$ square matrix A an element of R (called the determinant of A). Formally, we define the determinant function $\det : M_n(R) \rightarrow R$ by the formula

$$\det(A) = \sum_{f \in S_n} \text{sgn}(f) A(f(1), 1) \cdot A(f(2), 2) \cdots \cdot A(f(n), n) \quad (A \in M_n(R)). \quad (5.1)$$

(See Chapter 2 for the definitions of S_n and $\text{sgn}(f)$.) The sums and products on the right side of formula (5.1) are carried out in the ring R . Moreover, $\text{sgn}(f)$ is interpreted as $+1_R$ or -1_R , where 1_R is the multiplicative identity of R .

The true motivation for this mysterious formula comes from the theory of exterior powers in multilinear algebra, as we will see in §20.17. For now, we describe a way of visualizing the formula in terms of the geometric layout of the matrix A . A typical term in the determinant formula arises by choosing one entry in each column of the matrix A and multiplying these n entries together. If $f(i)$ is the row of the entry chosen in column i , for each $i \in [n]$, we obtain the product $\prod_{i=1}^n A(f(i), i)$. In order for this product to contribute to the determinant, f is required to be a permutation, which means that the entries picked in the n columns must lie in *distinct* rows. We attach the sign $\text{sgn}(f)$ to each term in the formula, where $\text{sgn}(f) = (-1)^{\text{inv}(f)}$ is the number of transpositions of adjacent elements needed to sort the list $[f(1), \dots, f(n)]$ into increasing order (see §2.6). The determinant of A is the sum of the signed terms coming from all $n!$ choices of the permutation f . We remark that we could also sum over all functions $f : [n] \rightarrow [n]$, keeping in mind the convention that $\text{sgn}(f) = 0$ when f is not a permutation.

For example, when $n = 3$,

$$\det \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} = aei + gbf + dhc - gec - dbi - ahf.$$

The term $-dbi$ comes from the permutation $f = [f(1), f(2), f(3)] = [2, 1, 3]$. The sign is $-1 = (-1)^1$, since it takes one transposition to sort the list $[2, 1, 3]$. The term gbf comes from the permutation $f = [3, 1, 2]$. The sign is $+1 = (-1)^2$, since it takes two transpositions of adjacent elements to sort the list $[3, 1, 2]$.

When $n = 4$, $\det(A)$ is a sum of $4! = 24$ terms. When $n = 5$, $\det(A)$ is a sum of $5! = 120$ terms. When $n = 8$, $\det(A)$ is a sum of $8! = 40320$ terms. We see that the explicit formula for the determinant quickly becomes unwieldy for doing computational work. However, it is very useful for theoretical purposes.

5.3 Diagonal and Triangular Matrices

If a matrix A has many zero entries, then many of the $n!$ terms in the formula for $\det(A)$ will be zero. This observation allows us to compute the determinants of certain special matrices very readily. For instance, consider the $n \times n$ identity matrix $I_n = I$. By definition,

$$\det(I) = \sum_{f \in S_n} \text{sgn}(f) I(f(1), 1) I(f(2), 2) \cdots I(f(n), n).$$

Since $I(i, j) = 0_R$ for all $i \neq j$ in $[n]$, the only way to get a nonzero summand is to choose f so that $f(1) = 1, f(2) = 2, \dots, f(n) = n$. In other words, f must be the identity permutation of S_n . Since $\text{sgn}(\text{id}) = +1$ and $I(i, i) = 1_R$ for all $i \in [n]$, we see that $\det(I) = 1_R$.

More generally, suppose D is a *diagonal* $n \times n$ matrix, which means that $D(i, j) = 0_R$ for all $i \neq j$ in $[n]$. The argument just used for I applies equally well to D , and shows that

$$\det(D) = D(1, 1)D(2, 2) \cdots D(n, n) = \prod_{i=1}^n D(i, i).$$

Thus, *the determinant of a diagonal matrix is the product of the entries on the diagonal.*

A slight adjustment of this technique allows us to compute the determinant of a *triangular* matrix. A matrix A is *upper-triangular* iff $A(i, j) = 0$ for all $i > j$ in $[n]$; it is *lower-triangular* iff $A(i, j) = 0$ for all $i < j$ in $[n]$. We will show that $\det(A)$ is the product of the diagonal entries of A in either case. Assume for instance that A is lower-triangular. For which $f \in S_n$ is the summand $A(f(1), 1)A(f(2), 2) \cdots A(f(n), n)$ in (5.1) nonzero? First, note that $A(f(n), n) \neq 0$ forces $f(n) \geq n$ and hence $f(n) = n$. Second, note that $A(f(n-1), n-1) \neq 0$ forces $f(n-1) \geq n-1$. Since f is a permutation and $f(n) = n$, we must have $f(n-1) = n-1$. Continuing by induction, fix $i \in [n]$ and assume that $f(j) = j$ for all $j > i$ in $[n]$. Then $A(f(i), i) \neq 0$ forces $f(i) \geq i$, and hence $f(i) = i$ since the values larger than i have already been used. After n steps, we see that $f = \text{id}$ is the only permutation that gives a nonzero contribution to $\det(A)$. Hence, $\det(A) = \prod_{i=1}^n A(i, i)$ as claimed. The argument for upper-triangular matrices is analogous.

5.4 Changing Variables

To prove many of the basic properties of determinants, we shall often need to use a “change of variable” to simplify certain sums of the form $\sum_{y \in S} H(y)$. Here, S is a finite set (usually S_n in the case of determinants), and $H : S \rightarrow R$ is a function that associates a summand $H(y)$ to each summation index $y \in S$. Suppose T is a set and $C : T \rightarrow S$ is a bijection (i.e., a function that is one-to-one and onto). Each $y \in S$ has the form $y = C(x)$ for a unique $x \in T$. Therefore, if a new summation variable x ranges over all elements of T exactly once,

then the values $y = C(x)$ will range over all elements of S exactly once. Using generalized associativity and commutativity of addition in the target ring R , it follows that

$$\sum_{y \in S} H(y) = \sum_{x \in T} H(C(x)). \quad (5.2)$$

In particular, if T and S happen to be the same set, we have

$$\sum_{x \in S} H(x) = \sum_{y \in S} H(y) = \sum_{x \in S} H(C(x)). \quad (5.3)$$

For example, let $S = S_n$, and let $C : S \rightarrow S$ be the inversion map given by $C(f) = f^{-1} \in S_n$ for every $f \in S_n$. We have $C(C(f)) = (f^{-1})^{-1} = f = \text{id}_{S_n}(f)$ for all $f \in S_n$. So $C \circ C = \text{id}_{S_n}$, which proves that C is a bijection on S_n and $C = C^{-1}$. Therefore,

$$\sum_{f \in S_n} H(f) = \sum_{f \in S_n} H(C(f)) = \sum_{f \in S_n} H(f^{-1}).$$

For another example, let $g \in S_n$ be fixed, and let $L_g : S_n \rightarrow S_n$ be the map given by $L_g(f) = g \circ f$ for all $f \in S_n$. The map L_g is called *left multiplication by g* . One may check that $L_{g^{-1}}$ (left multiplication by g^{-1}) is the two-sided inverse for L_g , so L_g is a bijection. The change-of-variable formula therefore gives

$$\sum_{f \in S_n} H(f) = \sum_{f \in S_n} H(L_g(f)) = \sum_{f \in S_n} H(g \circ f).$$

Exactly the same techniques apply to the simplification of finite products $\prod_{x \in S} H(x)$, assuming that the values $H(x)$ all lie in a given *commutative* ring R . Specifically, if S is finite and $C : S \rightarrow S$ is any bijection, we have

$$\prod_{x \in S} H(x) = \prod_{x \in S} H(C(x)). \quad (5.4)$$

5.5 Transposes and Determinants

In this section, we give an alternate formula for $\det(A)$ by applying an appropriate change of variable to the defining formula for the determinant. As a consequence, we will obtain the identity $\det(A^T) = \det(A)$.

Consider once again the definition (5.1) of the determinant of $A \in M_n(R)$. Writing the term $A(f(1), 1) \cdots A(f(n), n)$ as a product over $i \in [n]$, we see that

$$\det(A) = \sum_{f \in S_n} \text{sgn}(f) \prod_{i \in [n]} A(f(i), i).$$

Since $f \in S_n$, the map $C = f^{-1} : [n] \rightarrow [n]$ is a bijection of the index set of the product operation in this formula. Therefore, applying (5.4) with x replaced by i , S replaced by $[n]$, $H(x)$ replaced by $A(f(i), i)$, and C replaced by f^{-1} , we see that

$$\prod_{i \in [n]} A(f(i), i) = \prod_{i \in [n]} A(f(f^{-1}(i)), f^{-1}(i)) = \prod_{i \in [n]} A(i, f^{-1}(i))$$

for each $f \in S_n$. Putting this into the determinant formula gives

$$\det(A) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{i \in [n]} A(i, f^{-1}(i)).$$

We now perform a change of variable on the outer summation. Let $C : S_n \rightarrow S_n$ be the bijection $C(f) = f^{-1}$ for $f \in S_n$. Applying (5.3), we obtain

$$\begin{aligned} \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{i \in [n]} A(i, f^{-1}(i)) &= \sum_{f \in S_n} \operatorname{sgn}(f^{-1}) \prod_{i \in [n]} A(i, (f^{-1})^{-1}(i)) \\ &= \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{i \in [n]} A(i, f(i)). \end{aligned}$$

(We have used the facts that $\operatorname{sgn}(f) = \operatorname{sgn}(f^{-1})$ and $(f^{-1})^{-1} = f$.) Our final conclusion is that

$$\det(A) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{i \in [n]} A(i, f(i)). \quad (5.5)$$

Intuitively, this result means that we can compute the determinant by: selecting one element from each *row* of A , say the element in column $f(i)$ of row i for each $i \in [n]$; multiplying these chosen elements together; multiplying by $\operatorname{sgn}(f)$ (which is zero if f is not a permutation); and summing over all terms obtained by making such choices. In contrast, in our original discussion, we chose one element from each *column* of A .

We are now ready to prove that $\det(A^T) = \det(A)$ for all $A \in M_n(R)$. Using the original formula (5.1) for $\det(A^T)$, then the definition of A^T , and then the new formula (5.5) for $\det(A)$, we compute:

$$\begin{aligned} \det(A^T) &= \sum_{f \in S_n} \operatorname{sgn}(f) A^T(f(1), 1) A^T(f(2), 2) \cdots A^T(f(n), n) \\ &= \sum_{f \in S_n} \operatorname{sgn}(f) A(1, f(1)) A(2, f(2)) \cdots A(n, f(n)) \\ &= \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{i \in [n]} A(i, f(i)) = \det(A). \end{aligned}$$

To summarize our progress so far, we have shown:

$$\det(A) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{i=1}^n A(f(i), i) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{i=1}^n A(i, f(i)) = \det(A^T).$$

The symmetric roles of the rows and columns in these formulas will allow us to deduce many “column-oriented” properties of determinants from the corresponding “row-oriented” properties. Examples of this row-column symmetry will appear in the coming sections.

We conclude this section by deriving formulas for the determinant of the conjugate or conjugate-transpose of a complex square matrix. Let us temporarily write $c(z)$ to denote the complex conjugate \bar{z} of $z \in \mathbb{C}$, so that $c(x+iy) = x-iy$ for $x, y \in \mathbb{R}$. One may routinely verify that $c(z+z') = c(z) + c(z')$, $c(zz') = c(z)c(z')$, and $c(az) = ac(z)$ for $z, z' \in \mathbb{C}$ and $a \in \mathbb{R}$. The first two formulas extend by induction to finite sums and products. We therefore

have

$$\begin{aligned}
\det(\bar{A}) &= \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{i \in [n]} \bar{A}(f(i), i) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{i \in [n]} c(A(f(i), i)) \\
&= \sum_{f \in S_n} \operatorname{sgn}(f) c\left(\prod_{i \in [n]} A(f(i), i)\right) \\
&= c\left(\sum_{f \in S_n} \operatorname{sgn}(f) \prod_{i \in [n]} A(f(i), i)\right) \\
&= c(\det(A)) = \overline{\det(A)}.
\end{aligned}$$

In words, the determinant of the matrix conjugate of A is the conjugate of the determinant of A . Since the conjugate-transpose A^* is the matrix conjugate of A^T , we may also conclude that

$$\det(A^*) = \overline{\det(A^T)} = \overline{\det(A)}.$$

5.6 Multilinearity and the Alternating Property

In this section and the next one, we investigate the relationship between the determinant of a square matrix A and the rows of A . Given $A \in M_n(R)$ and $i \in [n]$, we let $A_{[i]} \in R^n$ be the i 'th row of A , so that $A_{[i]}(j) = A(i, j)$ for all $j \in [n]$. In this section, we will often identify the matrix A with the ordered list of rows $(A_{[1]}, A_{[2]}, \dots, A_{[n]}) \in (R^n)^n$. Thus, the determinant map can be regarded as a function $\det : (R^n)^n \rightarrow R$ of n variables, where each “variable” $A_{[i]} = A_i$ is an element of R^n that encodes the i 'th row of A . Using the row notation, our formulas (5.1) and (5.5) for the determinant of A become

$$\det(A_1, \dots, A_n) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{k=1}^n A_{f(k)}(k) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{k=1}^n A_k(f(k)).$$

We will now show that the determinant map is *R-multilinear* in its n arguments A_1, \dots, A_n . multilinearity means that, if an index $i \in [n]$ and the rows A_j for all $j \neq i$ are held fixed, then the determinant function is *R-linear* with respect to the remaining argument A_i , which is allowed to vary over R^n . More precisely, fix $i \in [n]$ and n -tuples $A_j \in R^n$ for all $j \neq i$ in $[n]$, and define $D : R^n \rightarrow R$ by $D(v) = \det(A_1, \dots, A_{i-1}, v, A_{i+1}, \dots, A_n)$ for $v \in R^n$; then we are asserting that D is an *R-linear* map (as defined in §1.7). To prove this, replace A_i by $v = (v(1), \dots, v(n))$ in the second formula for $\det(A_1, \dots, A_n)$ above, obtaining

$$D(v) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{k=1}^{i-1} A_k(f(k)) v(f(i)) \prod_{k=i+1}^n A_k(f(k)) = \sum_{f \in S_n} v(f(i)) \operatorname{sgn}(f) \prod_{\substack{k \in [n] \\ k \neq i}} A_k(f(k)).$$

Each summand here consists of one of the components of v multiplied by a scalar in R . Grouping together summands that involve the same component of v , we see that we can write $D(v) = \sum_{j=1}^n c_j v(j)$, where each c_j is a scalar in R that depends on the fixed A_k 's

with $k \neq i$, but not on v . To be explicit, we have (for each $j \in [n]$)

$$c_j = \sum_{\substack{f \in S_n: \\ f(i)=j}} \operatorname{sgn}(f) \prod_{\substack{k \in [n] \\ k \neq i}} A_k(f(k)).$$

The R -linearity of the function D now follows from the computations:

$$\begin{aligned} D(v+w) &= \sum_{j=1}^n c_j(v+w)(j) = \sum_{j=1}^n c_j v(j) + \sum_{j=1}^n c_j w(j) = D(v) + D(w) \quad (v, w \in R^n); \\ D(av) &= \sum_{j=1}^n c_j(av)(j) = a \sum_{j=1}^n c_j v(j) = aD(v) \quad (v \in R^n, a \in R). \end{aligned}$$

Since D is linear, $D(0_{R^n}) = 0_R$, which means that *a matrix containing a row of zeroes has determinant zero*.

Here is an example of the formula $D(v) = \sum_{j=1}^n c_j v(j)$ in the case where $A \in M_3(\mathbb{R})$ and $i = 3$. Fix rows $A_1 = (4, 1, -3)$ and $A_2 = (0, 2, 5)$; then for $v = (v_1, v_2, v_3) \in \mathbb{R}^3$,

$$D(v) = \det \begin{bmatrix} 4 & 1 & -3 \\ 0 & 2 & 5 \\ v_1 & v_2 & v_3 \end{bmatrix} = 11v_1 - 20v_2 + 8v_3,$$

and this is an \mathbb{R} -linear function of the row vector v .

Another often-used fact about determinants is the *alternating property*: if $A_i = A_j$ for some $i \neq j$ in $[n]$, then $\det(A_1, \dots, A_n) = 0$. In other words, *the determinant of a matrix with two equal rows is zero*. Under the assumption $A_i = A_j$, we can write the formula for $\det(A_1, \dots, A_n)$ as

$$\sum_{f \in S_n} \operatorname{sgn}(f) \prod_{k \in [n]} A_k(f(k)) = \sum_{f \in S_n} \operatorname{sgn}(f) A_i(f(i)) A_i(f(j)) \prod_{\substack{k \in [n] \\ k \neq i, j}} A_k(f(k)).$$

Let us compare the summand coming from some $f \in S_n$ with the summand coming from $f' = f \circ (i, j)$. We have $f'(i) = f(j)$, $f'(j) = f(i)$, $f'(k) = f(k)$ for $k \neq i, j$, and $\operatorname{sgn}(f') = -\operatorname{sgn}(f)$ (see §2.7). It follows that

$$\operatorname{sgn}(f') A_i(f'(i)) A_i(f'(j)) \prod_{\substack{k \in [n] \\ k \neq i, j}} A_k(f'(k)) = -\operatorname{sgn}(f) A_i(f(j)) A_i(f(i)) \prod_{\substack{k \in [n] \\ k \neq i, j}} A_k(f(k)).$$

Since R is commutative, the summands for f and f' cancel each other. Pairing off all $n!$ terms in the sum for $\det(A)$ in this way, we see that all terms cancel, hence $\det(A) = 0$.

5.7 Elementary Row Operations and Determinants

Now that we know the determinant is alternating and multilinear in the rows of A , we can deduce some well-known results concerning the effect of elementary row operations on determinants. The first type of elementary row operation multiplies the i 'th row of A by an invertible scalar $c \in R$. From the multilinearity property, we know

$$\det(A_1, \dots, cA_i, \dots, A_n) = c \det(A_1, \dots, A_i, \dots, A_n)$$

for any scalar $c \in R$ (invertible or not). In words, multiplying *one row* of a matrix by a scalar multiplies the determinant of the matrix by that scalar. Multiplying the *whole matrix* by a given scalar amounts to multiplying each of the n rows by that scalar. Iteration of the result just mentioned therefore gives

$$\det(cA) = c^n \det(A) \quad (A \in M_n(R)).$$

The second type of elementary row operation interchanges two distinct rows of A , say row i and row j . We claim this operation causes the determinant of A to switch sign, i.e.,

$$\det(A_1, \dots, A_i, \dots, A_j, \dots, A_n) = -\det(A_1, \dots, A_j, \dots, A_i, \dots, A_n) \in R.$$

To prove this, fix $i \neq j$ in $[n]$, fix the other rows A_k for $k \neq i, j$ in $[n]$, and define $D : R^n \times R^n \rightarrow R$ by $D(v, w) = \det(A_1, \dots, v, \dots, w, \dots, A_n)$ for $v, w \in R^n$, where v and w occur in position i and position j , respectively. By the alternating property already proved, $D(v, v) = 0$ for all $v \in R^n$. Furthermore, D is linear in each of its arguments, since the determinant is multilinear. We therefore have

$$\begin{aligned} 0 &= D(A_i + A_j, A_i + A_j) = D(A_i + A_j, A_i) + D(A_i + A_j, A_j) \\ &= D(A_i, A_i) + D(A_j, A_i) + D(A_i, A_j) + D(A_j, A_j) \\ &= 0 + D(A_j, A_i) + D(A_i, A_j) + 0. \end{aligned}$$

Rearranging, we get $D(A_i, A_j) = -D(A_j, A_i)$, and the result we want follows from the definition of D .

The third type of elementary row operation adds c times row j to row i of a matrix, for $i \neq j$ in $[n]$ and any $c \in R$. We claim this operation does not affect the determinant of the matrix, i.e.,

$$\det(A_1, \dots, A_i + cA_j, \dots, A_j, \dots, A_n) = \det(A_1, \dots, A_i, \dots, A_j, \dots, A_n) \in R.$$

This fact follows from multilinearity and the alternating property. For, letting D be as in the last paragraph, we have

$$D(A_i + cA_j, A_j) = D(A_i, A_j) + cD(A_j, A_j) = D(A_i, A_j) + c \cdot 0 = D(A_i, A_j),$$

as needed.

In closing, let us interpret these results in terms of elementary matrices. Recall from §4.9 that if A' results from $A \in M_n(R)$ by performing some elementary row operation, then $A' = EA$ where E is the *elementary matrix* obtained from I_n by applying the same row operation that took A to A' . We wish to show that $\det(A') = \det(E) \det(A)$ for all three types of elementary operations. If A' was obtained by multiplying row i of A by $c \in R$ (where c is invertible), then E is a diagonal matrix with a c in the i, i -position and 1's elsewhere on the diagonal. We have proved in this section that $\det(A') = c \det(A)$, whereas $\det(E) = c$ by the formula for the determinant of a diagonal matrix. Therefore $\det(A') = \det(E) \det(A)$ in this case. If, instead, A' was obtained from A by interchanging row i and row j (where $i \neq j$ in $[n]$), we have seen that $\det(A') = -\det(A)$. On the other hand, since the corresponding elementary matrix E is obtained from I_n by interchanging two rows, we have $\det(E) = -\det(I_n) = -1_R$. Therefore, $\det(A') = \det(E) \det(A)$ in this case as well. If, finally, A' was obtained from A by adding c times row j to row i (for some $i \neq j$ in $[n]$ and $c \in R$), we have seen that $\det(A') = \det(A)$. In this case, E is obtained from I_n by the same type of operation, and so $\det(E) = \det(I_n) = 1_R$. Therefore, $\det(A') = \det(E) \det(A)$ in this third case. In summary, *for all square matrices A and all elementary matrices E of the same size as A , we have*

$$\det(EA) = \det(E) \det(A).$$

When R is a field, we can use this fact to prove the general product formula $\det(BA) = \det(B)\det(A)$ for all $A, B \in M_n(R)$ (see §5.9).

5.8 Determinant Properties Involving Columns

The discussion in the last two sections, relating determinants to the rows of a matrix, applies equally well to the columns of a matrix. Given $A \in M_n(R)$ and $j \in [n]$, we let $A^{[j]} \in R^n$ be the j 'th column of A , so that $A^{[j]}(i) = A(i, j)$ for all $i \in [n]$. We now identify the matrix A with the ordered list of its columns, writing $A = (A^{[1]}, \dots, A^{[n]}) \in (R^n)^n$. Using this column notation, our formulas (5.1) and (5.5) for the determinant become

$$\det(A) = \det(A^{[1]}, \dots, A^{[n]}) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{k=1}^n A^{[k]}(f(k)) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{k=1}^n A^{[f(k)]}(k).$$

Starting from these formulas, we can now repeat all the proofs in the last two sections verbatim, merely changing all subscripts (row indices) to superscripts (column indices).

We thereby conclude that: the determinant is an R -multilinear function of the n columns of A ; $\det(A) = 0$ if any two columns of A are equal or if A has a column of zeroes; multiplying a column of A by $c \in R$ multiplies the corresponding determinant by c ; interchanging two columns of A flips the sign of the determinant; and adding a scalar multiple of one column to a different column does not change the determinant. Since elementary column operations correspond to right multiplication by elementary matrices (see §4.9), the last three observations imply (as above) that $\det(AE) = \det(A)\det(E)$ for all matrices $A \in M_n(R)$ and all *elementary* matrices $E \in M_n(R)$.

All of these results can also be derived quickly from the corresponding results for rows by invoking the identity $\det(A) = \det(A^T)$. For instance, if A has two equal columns, then A^T has two equal rows, and so $\det(A) = \det(A^T) = 0$. For another example of this proof method, assume A is arbitrary and E is elementary. Then E^T is also elementary (as is readily checked), so that

$$\det(AE) = \det((AE)^T) = \det(E^T A^T) = \det(E^T) \det(A^T) = \det(A) \det(E).$$

This computation relies heavily on the commutativity of the ring R (cf. Exercise 56).

5.9 Product Formula via Elementary Matrices

In this section only, we consider matrices with entries in a *field* F . We will use facts about elementary matrices and row-reduction to prove the *product formula* $\det(AB) = \det(A)\det(B)$ for $A, B \in M_n(F)$, and to prove that A is invertible iff $\det(A) \neq 0_F$. Proofs of corresponding results for matrices with entries in any commutative ring will be given later.

First, we indicate a method for computing determinants of matrices in $M_n(F)$ that is, in general, much more efficient than computing all $n!$ terms in the defining formula. Given a square matrix A , perform the *Gaussian elimination algorithm*, using a finite sequence of elementary row operations to bring A into echelon form. Since A is square, the echelon form must be an upper-triangular matrix. We can use the results in §5.7 to keep track of how

the determinant of the matrix changes as each row operation is applied. As we have seen, for every elementary row operation, the determinant is multiplied by some *nonzero* scalar in the field F (possibly the scalar 1). When a triangular matrix is reached, we can find the determinant quickly by multiplying the diagonal entries.

For example, let the matrix $A = \begin{bmatrix} 0 & 2 & 1 \\ 3 & 6 & -3 \\ 2 & 5 & 1 \end{bmatrix}$ have unknown determinant $\det(A) = d$.

We perform Gaussian elimination steps on A as shown below, keeping track of the effect of each elementary row operation on the determinant:

$$\begin{array}{ccc} \left[\begin{array}{ccc} 0 & 2 & 1 \\ 3 & 6 & -3 \\ 2 & 5 & 1 \end{array} \right] & \xrightarrow{\text{R1} \leftrightarrow \text{R2}} & \left[\begin{array}{ccc} 3 & 6 & -3 \\ 0 & 2 & 1 \\ 2 & 5 & 1 \end{array} \right] \\ \text{det} = d & & \text{det} = -d \end{array} \quad \begin{array}{ccc} & \xrightarrow{\text{R1} \times 1/3} & \left[\begin{array}{ccc} 1 & 2 & -1 \\ 0 & 2 & 1 \\ 2 & 5 & 1 \end{array} \right] \\ & & \text{det} = -d/3 \end{array}$$

$$\begin{array}{ccc} \xrightarrow{\text{R3} \rightarrow 2\text{R1}} \left[\begin{array}{ccc} 1 & 2 & -1 \\ 0 & 2 & 1 \\ 0 & 1 & 3 \end{array} \right] & \xrightarrow{\text{R2} \leftrightarrow \text{R3}} & \left[\begin{array}{ccc} 1 & 2 & -1 \\ 0 & 1 & 3 \\ 0 & 2 & 1 \end{array} \right] \\ \text{det} = -d/3 & & \text{det} = d/3 \end{array} \quad \begin{array}{ccc} & \xrightarrow{\text{R3} \rightarrow 2\text{R2}} & \left[\begin{array}{ccc} 1 & 2 & -1 \\ 0 & 1 & 3 \\ 0 & 0 & -5 \end{array} \right] \\ & & \text{det} = d/3 \end{array}$$

The final matrix is triangular, with determinant $1 \cdot 1 \cdot (-5) = -5$. Since this also equals $d/3$, we see that $\det(A) = -15$.

Returning to the discussion of a general matrix $A \in M_n(F)$, we can continue to perform elementary row operations to reach the reduced row-echelon form A_{ech} , defined in §4.10. We know $A_{\text{ech}} = PA$ for some invertible matrix $P \in M_n(F)$, where P is the product of the elementary matrices associated with the elementary row operations used to reduce A . So, A is invertible iff A_{ech} is invertible. Consequently, if A is invertible, then A_{ech} cannot have any zero rows. Since A is square, we must have $A_{\text{ech}} = I_n$, which has determinant 1_F . The determinant of A differs from this by some factor that is a product of nonzero elements of F ; hence, A invertible implies $\det(A) \neq 0_F$. Conversely, suppose A is not invertible. Then A_{ech} is not invertible, hence must be an upper-triangular matrix with at least one zero on the main diagonal, so $\det(A_{\text{ech}}) = 0_F$. Again, the determinant of A is a scalar multiple of this determinant, so A non-invertible implies $\det(A) = 0_F$.

We now turn to the proof of the product formula $\det(AB) = \det(A)\det(B)$ for $A, B \in M_n(F)$. Recall that the product formula has already been proved when A is an *elementary* matrix and B is arbitrary (§5.7). Next we show that the formula holds when $A = E_1 \cdots E_k$ is a finite product of elementary matrices, and B is arbitrary. We prove this by induction on k , the case $k = 1$ already being known. Assuming the result holds for $k - 1$ factors, and using associativity of matrix multiplication and the result for $k = 1$, we compute

$$\begin{aligned} \det(AB) &= \det(E_1(E_2 \cdots E_k B)) = \det(E_1) \det((E_2 \cdots E_k)B) \\ &= \det(E_1) \det(E_2 \cdots E_k) \det(B) = \det(E_1(E_2 \cdots E_k)) \det(B) \\ &= \det(A) \det(B). \end{aligned}$$

Using row-reduction as above, one can show that $A \in M_n(R)$ is invertible iff A is a finite product of elementary matrices (see §4.11). Therefore, we have proved that the product formula $\det(AB) = \det(A)\det(B)$ holds when A is invertible and B is arbitrary. Similarly, the formula holds when B is invertible and A is arbitrary. Finally, consider the case where neither A nor B is invertible. As B is not invertible, there exists some nonzero $x \in F^n$ with $Bx = 0$ (see §4.13). Hence $(AB)x = 0$, so AB is not invertible. It follows that $\det(AB) = 0_F = 0 \cdot 0 = \det(A)\det(B)$.

5.10 Laplace Expansions

We are now ready to prove the *Laplace expansions*, which are often used as the definition of the determinant in introductory treatments. Given $A \in M_n(R)$, let $A[i|j]$ be the matrix in $M_{n-1}(R)$ obtained by deleting row i and column j of A . Formally, the u, v -entry of $A[i|j]$ is $A(u, v)$ if $u < i$ and $v < j$; $A(u+1, v)$ if $i \leq u \leq n-1$ and $v < j$; $A(u, v+1)$ if $u < i$ and $j \leq v \leq n-1$; and $A(u+1, v+1)$ if $i \leq u \leq n-1$ and $j \leq v \leq n-1$. Writing $\chi(P) = 1$ if P is true and $\chi(P) = 0$ if P is false (for any logical statement P), we can define $A[i|j]$ more succinctly by the formula

$$A[i|j](u, v) = A(u + \chi(u \geq i), v + \chi(v \geq j)) \quad (u, v \in [n-1]).$$

With this notation in place, we have, for each $i \in [n]$, the *Laplace expansion along row i* :

$$\det(A) = \sum_{j=1}^n A(i, j)(-1)^{i+j} \det(A[i|j]).$$

We also have, for each $j \in [n]$, the *Laplace expansion along column j* :

$$\det(A) = \sum_{i=1}^n A(i, j)(-1)^{i+j} \det(A[i|j]).$$

Intuitively, these formulas arise by grouping the $n!$ terms in the definition of $\det(A)$ based on which entry we choose to use from row i (or column j). We prove these formulas in several steps.

Step 1: We prove the row expansion formula for $i = n$. Starting with the defining formula

$$\det(A) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{k \in [n]} A(f(k), k),$$

we can group together the terms in this sum that involve $A(n, j)$ for each j between 1 and n . The entry $A(n, j)$ appears in the term indexed by $f \in S_n$ iff $f(j) = n$. Therefore, $\det(A) = \sum_{j=1}^n A(n, j)U_j$, where

$$U_j = \sum_{\substack{f \in S_n: \\ f(j)=n}} \operatorname{sgn}(f) \prod_{\substack{k \in [n]: \\ k \neq j}} A(f(k), k).$$

Step 1 will be completed once we show that this expression for U_j equals $(-1)^{n+j} \det(A[n|j])$.

From now on, keep j fixed. Introduce the notation $S = \{f \in S_n : f(j) = n\}$ and $T = S_{n-1}$. Define a bijection $C : T \rightarrow S$ as follows. View an element $g \in S_{n-1}$ as a list $[g(1), \dots, g(n-1)]$ that uses the numbers 1 through $n-1$ once each. Define $C(g)$ to be the list

$$g' = [g(1), \dots, g(j-1), n, g(j), \dots, g(n-1)],$$

which is an element of S . C is a bijection; the inverse map deletes the j 'th element from a list $f \in S$, which must be the integer n .

Applying the change-of-variable formula (5.2), and writing g' for $C(g)$, we see that

$$U_j = \sum_{f \in S} \operatorname{sgn}(f) \prod_{\substack{k \in [n] \\ k \neq j}} A(f(k), k) = \sum_{g \in T} \operatorname{sgn}(g') \prod_{\substack{k \in [n] \\ k \neq j}} A(g'(k), k). \quad (5.6)$$

To continue simplifying, let us first relate $\text{sgn}(g')$ to $\text{sgn}(g)$. Recall that $\text{sgn}(g) = (-1)^{\text{inv}(g)}$, where $\text{inv}(g)$ is the minimum number of basic transposition moves needed to sort the list $g = [g(1), \dots, g(n-1)]$ (see §2.6). Similarly, $\text{sgn}(g') = (-1)^{\text{inv}(g')}$, where $\text{inv}(g')$ is the minimum number of basic transposition moves needed to sort the list

$$g' = [g(1), \dots, g(j-1), n, g(j), \dots, g(n-1)].$$

The latter list may be sorted by first moving the largest entry n from position j to position n (which requires $n-j$ basic transposition moves) and then sorting the first $n-1$ elements of the resulting list. Since the resulting list is $[g(1), \dots, g(n-1), n]$, the second stage of the sorting requires $\text{inv}(g)$ moves. Therefore, $\text{inv}(g') = \text{inv}(g) + n - j$. Raising -1_R to this power and noting that $(-1)^{n-j} = (-1)^{n+j}$ in any ring, we see that

$$\text{sgn}(g') = (-1)^{n+j} \text{sgn}(g). \quad (5.7)$$

Next, what is $A(g'(k), k)$ in terms of g and $A[n|j]$? If $1 \leq k < j$, then $g'(k) = g(k)$ by definition. Also $g(k) < n$, so $A[n|j](g(k), k) = A(g(k), k) = A(g'(k), k)$ in this case. On the other hand, if $j < k \leq n$, then $g'(k) = g(k-1) < n$ and $A[n|j](g(k-1), k-1) = A(g(k-1), k) = A(g'(k), k)$ by definition. We can therefore write

$$\begin{aligned} \prod_{\substack{k \in [n] \\ k \neq j}} A(g'(k), k) &= \prod_{1 \leq k < j} A[n|j](g(k), k) \prod_{j < k \leq n} A[n|j](g(k-1), k-1) \\ &= \prod_{1 \leq k < j} A[n|j](g(k), k) \prod_{j \leq k \leq n-1} A[n|j](g(k), k) \\ &= \prod_{k \in [n-1]} A[n|j](g(k), k). \end{aligned}$$

Using this and (5.7) in (5.6), we see that

$$U_j = (-1)^{n+j} \sum_{g \in S_{n-1}} \text{sgn}(g) \prod_{k \in [n-1]} A[n|j](g(k), k) = (-1)^{n+j} \det(A[n|j]).$$

Step 2: We prove that if the Laplace expansion formula holds for row $i > 1$, then it holds for row $i-1$. Combining this with step 1, we can conclude that the Laplace expansion is valid for any row of A . To prove step 2, assume that $1 < i \leq n$ is fixed and that

$$\det(B) = \sum_{j=1}^n B(i, j)(-1)^{i+j} \det(B[i|j])$$

is known to be true for any $B \in M_n(R)$. Given A , let B be the matrix obtained from A by interchanging row $i-1$ and row i . We know that $\det(B) = -\det(A)$. Moreover, it follows from the definitions that $B[i|j] = A[i-1|j]$ and $B(i, j) = A(i-1, j)$ for all $j \in [n]$. Therefore,

$$\begin{aligned} \det(A) &= -\det(B) = \sum_{j=1}^n B(i, j)(-1)^{i+j+1} \det(B[i|j]) \\ &= \sum_{j=1}^n A(i-1, j)(-1)^{(i-1)+j} \det(A[i-1|j]). \end{aligned}$$

This is the Laplace expansion formula for A along row $i-1$.

Step 3: We prove the Laplace expansion formula along column k , for each $k \in [n]$. From steps 1 and 2, we already know that

$$\det(B) = \sum_{s=1}^n B(k, s)(-1)^{k+s} \det(B[k|s])$$

holds for any B and any row index k . Given A , take $B = A^T$. The definitions show that $B(k, s) = A(s, k)$ and $B[k|s] = A[s|k]^T$ (since deleting row k and column s of $B = A^T$ has the same effect as deleting row s and column k of A , and then transposing). We now compute

$$\begin{aligned} \det(A) &= \det(B) = \sum_{s=1}^n B(k, s)(-1)^{k+s} \det(B[k|s]) \\ &= \sum_{s=1}^n A(s, k)(-1)^{s+k} \det(A[s|k]^T) = \sum_{i=1}^n A(i, k)(-1)^{i+k} \det(A[i|k]). \end{aligned}$$

This is the Laplace expansion formula for A along column k .

5.11 Classical Adjoints and Inverses

Given a matrix $A \in M_n(R)$, we now define another $n \times n$ matrix called the *classical adjoint* or *adjunct* of A and denoted $\text{adj}(A)$ or A' . For $i, j \in [n]$, we set $A'(i, j) = (-1)^{i+j} \det(A[j|i])$, which is the determinant of the matrix obtained by deleting row j and column i of A , times a certain sign. This matrix is “almost” the inverse of A ; more precisely, we will show that $AA' = A'A = \det(A)I_n \in M_n(R)$.

We prove that $AA' = \det(A)I_n$ in two stages. First, consider the i, j -entry of AA' for some $i \neq j$ in $[n]$. The definition of matrix multiplication gives

$$(AA')(i, j) = \sum_{k=1}^n A(i, k)A'(k, j) = \sum_{k=1}^n A(i, k)(-1)^{j+k} \det(A[j|k]).$$

Let C be the matrix obtained from A by replacing row j of A by row i . We see that $C[j|k] = A[j|k]$ and $C(j, k) = A(i, k)$ for all $k \in [n]$. Since C has two equal rows, $\det(C) = 0$. On the other hand, the Laplace expansion along row j of C gives

$$0 = \det(C) = \sum_{k=1}^n C(j, k)(-1)^{j+k} \det(C[j|k]) = \sum_{k=1}^n A(i, k)(-1)^{j+k} \det(A[j|k]).$$

Comparing to the previous formula, we see that $(AA')(i, j) = 0 = (\det(A)I_n)(i, j)$.

The second stage is to compute $(AA')(i, i)$, where $i \in [n]$. Here we find that

$$(AA')(i, i) = \sum_{k=1}^n A(i, k)A'(k, i) = \sum_{k=1}^n A(i, k)(-1)^{i+k} \det(A[i|k]) = \det(A),$$

where the last equality is Laplace expansion along row i . Thus, $(AA')(i, i) = \det(A) = (\det(A)I_n)(i, i)$, as needed.

The companion formula $A'A = I_n$ can be proved similarly, using Laplace expansions

down columns. Alternatively, the definitions readily give $(A')^T = (A^T)'$, so that the result already proved (applied to A^T) gives

$$(A'A)^T = A^T(A')^T = A^T(A^T)' = \det(A^T)I_n = \det(A)I_n.$$

Transposing both sides gives $A'A = \det(A)I_n^T = \det(A)I_n$.

We now discuss some consequences of the formulas $A'A = \det(A)I_n = A'A$. First, if $\det(A)$ is an *invertible* element in the ring R , then A is an invertible matrix. Indeed, multiplying the previous formulas by the inverse of $\det(A)$ in R , we see that the inverse of A is given explicitly by $A^{-1} = \det(A)^{-1}A'$, i.e.,

$$A^{-1}(i,j) = (-1)^{i+j} \det(A[j|i]) \det(A)^{-1} \quad (i,j \in [n]).$$

Conversely, if A is invertible, then $\det(A)$ is invertible in R . The converse follows readily using the product formula $\det(AB) = \det(A)\det(B)$ (proved in §5.13 below): for if $AB = I_n = BA$, taking determinants yields $\det(A)\det(B) = \det(I_n) = \det(B)\det(A)$. Since $\det(I_n) = 1_R$, we see that $\det(B)$ is a two-sided inverse for $\det(A)$ in the ring R . In summary, *for any commutative ring R , the matrix A is invertible in $M_n(R)$ if and only if $\det(A)$ is invertible in R .*

For example, consider the matrix $A = \begin{bmatrix} 0 & 2 & 1 \\ 3 & 6 & -3 \\ 2 & 5 & 1 \end{bmatrix} \in M_3(\mathbb{Z})$. The classical adjoint of A is $A' = \begin{bmatrix} 21 & 3 & -12 \\ -9 & -2 & 3 \\ 3 & 4 & -6 \end{bmatrix}$; for instance, $A'(2,3) = (-1)^{2+3} \det \begin{bmatrix} 0 & 1 \\ 3 & -3 \end{bmatrix} = 3$. We found earlier that $\det(A) = -15$, and one readily checks that $AA' = A'A = -15I_3$. Since -15 is not invertible in the ring \mathbb{Z} , A has no inverse in the ring $M_3(\mathbb{Z})$. On the other hand, A is invertible in the ring $M_3(\mathbb{R})$, with inverse $A^{-1} = (-1/15)A'$.

5.12 Cramer's Rule

One popular application of the explicit formula for A^{-1} is *Cramer's rule* for solving a nonsingular system of n linear equations in n unknowns (*nonsingular* means the system has a unique solution). Write this system in matrix notation as $A\vec{x} = \vec{b}$, where $A \in M_n(R)$ has invertible determinant, $\vec{x} \in R^n$ is a column vector of “unknowns,” and $\vec{b} \in R^n$ is a given column vector. To solve for a particular unknown x_i , multiply both sides by A^{-1} on the left, obtaining $\vec{x} = A^{-1}\vec{b}$, and take the i 'th component:

$$x_i = (A^{-1}\vec{b})_i = \sum_{k=1}^n A^{-1}(i,k)b_k.$$

Using the explicit formula for A^{-1} , we get

$$x_i = \det(A)^{-1} \sum_{k=1}^n b_k (-1)^{i+k} \det(A[k|i]).$$

To obtain Cramer's rule, let C_i be the matrix whose columns are

$$(A^{[1]}, \dots, A^{[i-1]}, \vec{b}, A^{[i+1]}, \dots, A^{[n]}).$$

Evidently, $C_i[k|i] = A[k|i]$ and $C_i(k, i) = b_k$ for all $k \in [n]$. Laplace expansion down column i of C_i therefore gives

$$\det(C_i) = \sum_{k=1}^n C_i(k, i)(-1)^{k+i} \det(C_i[k|i]) = \sum_{k=1}^n b_k (-1)^{i+k} \det(A[k|i]).$$

Comparing this to the formula for x_i , we finally get $x_i = \det(A)^{-1} \det(C_i)$ for $i \in [n]$. In words, x_i may be computed by replacing the i 'th column of A by the coefficient vector \vec{b} , taking the determinant of this matrix, and dividing that by the determinant of A itself. We stress that Cramer's rule only applies when $\det(A)$ is invertible in R . Moreover, when R is a field, Gaussian elimination is often a more efficient way to solve the linear system compared to the multiple determinant evaluations required by Cramer's rule.

As an example of Cramer's rule, let us solve the linear system of equations

$$\begin{cases} 2x_2 & +x_3 = 9 \\ 3x_1 & +6x_2 -3x_3 = -2 \\ 2x_1 & +5x_2 +x_3 = 4. \end{cases}$$

This system has the form $A\vec{x} = \vec{b}$, where $A = \begin{bmatrix} 0 & 2 & 1 \\ 3 & 6 & -3 \\ 2 & 5 & 1 \end{bmatrix}$ and $\vec{b} = \begin{bmatrix} 9 \\ -2 \\ 4 \end{bmatrix}$. Earlier, we calculated $\det(A) = -15$. Hence, Cramer's rule gives

$$\begin{aligned} x_1 &= (-1/15) \det \begin{bmatrix} 9 & 2 & 1 \\ -2 & 6 & -3 \\ 4 & 5 & 1 \end{bmatrix} = \frac{135}{-15} = -9; \\ x_2 &= (-1/15) \det \begin{bmatrix} 0 & 9 & 1 \\ 3 & -2 & -3 \\ 2 & 4 & 1 \end{bmatrix} = \frac{-65}{-15} = \frac{13}{3}; \\ x_3 &= (-1/15) \det \begin{bmatrix} 0 & 2 & 9 \\ 3 & 6 & -2 \\ 2 & 5 & 4 \end{bmatrix} = \frac{-5}{-15} = \frac{1}{3}. \end{aligned}$$

5.13 Product Formula for Determinants

Let $A, B \in M_n(R)$, where R is any commutative ring. We now give a direct computational proof of the product formula $\det(AB) = \det(A)\det(B)$. We start with the definition (5.1):

$$\det(AB) = \sum_{f \in S_n} \operatorname{sgn}(f)(AB)(f(1), 1)(AB)(f(2), 2) \cdots (AB)(f(n), n).$$

Next, use the definition of matrix multiplication to write each factor as

$$(AB)(f(j), j) = \sum_{k_j=1}^n A(f(j), k_j)B(k_j, j).$$

Putting these expressions into the product above and using the distributive law repeatedly (see Exercise 5), we get:

$$\det(AB) = \sum_{f \in S_n} \sum_{k_1=1}^n \cdots \sum_{k_n=1}^n \operatorname{sgn}(f) A(f(1), k_1) B(k_1, 1) \cdots A(f(n), k_n) B(k_n, n).$$

We now reorder the sums and products, using commutativity of addition and multiplication in R :

$$\det(AB) = \sum_{k_1=1}^n \cdots \sum_{k_n=1}^n \sum_{f \in S_n} \operatorname{sgn}(f) A(f(1), k_1) \cdots A(f(n), k_n) B(k_1, 1) \cdots B(k_n, n).$$

Note that we can regard the n outer summations over integers $k_j \in [n]$ as a single summation over all lists $k = [k_1, k_2, \dots, k_n] \in [n]^n$. Also, we can pull out the factors $B(k_j, j)$ that do not depend on the inner summation index f . We get:

$$\det(AB) = \sum_{k \in [n]^n} B(k_1, 1) \cdots B(k_n, n) \left[\sum_{f \in S_n} \operatorname{sgn}(f) A(f(1), k_1) \cdots A(f(n), k_n) \right].$$

Given a fixed list $k \in [n]^n$, note that the term in brackets is the determinant of the matrix C whose columns are $(A^{[k_1]}, \dots, A^{[k_n]})$. To see why, note that $C(i, j) = A(i, k_j)$ for $i, j \in [n]$, and so

$$\det(C) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{j \in [n]} C(f(j), j) = \sum_{f \in S_n} \operatorname{sgn}(f) \prod_{j \in [n]} A(f(j), k_j).$$

Now, if $k_i = k_j$ for any $i \neq j$, we know that $\det(C) = 0$ because it has two equal columns. Therefore, we can throw away all terms in the outer summation indexed by lists $k = [k_1, \dots, k_n]$ with a repeated entry. The surviving summands are indexed by permutations $k = [k_1, \dots, k_n] \in S_n$. Suppose we sort the list $[k_1, \dots, k_n]$ into increasing order by interchanging adjacent elements, and at the same time interchanging the corresponding adjacent columns in C . Each such interchange will multiply $\det(C)$ by -1 . We need $\operatorname{inv}(k)$ interchanges to reach the list $[1, 2, \dots, n]$ (see §2.6), and at that point we will have transformed C back into the matrix A whose columns are in their original order. Therefore, $\det(C) = (-1)^{\operatorname{inv}(k)} \det(A) = \operatorname{sgn}(k) \det(A)$. Replacing the term in square brackets by this formula, and factoring $\det(A)$ out of the main sum, we now see that

$$\det(AB) = \det(A) \sum_{k \in S_n} \operatorname{sgn}(k) B(k_1, 1) \cdots B(k_n, n).$$

The remaining sum is precisely the definition of $\det(B)$, so we are done.

A very short but abstract proof of the product formula will be given in §20.17.

5.14 Cauchy–Binet Formula

The *Cauchy–Binet formula* is a generalization of the product formula $\det(AB) = \det(A) \det(B)$ to products of rectangular matrices. Suppose $m \leq n$, A is an $m \times n$ matrix, and B is an $n \times m$ matrix. Let J be the set of all strictly increasing lists $j = [j_1, j_2, \dots, j_m]$ with entries in $[n]$. The Cauchy–Binet formula is

$$\det(AB) = \sum_{j \in J} \det(A^{[j_1]}, A^{[j_2]}, \dots, A^{[j_m]}) \det(B_{[j_1]}, B_{[j_2]}, \dots, B_{[j_m]}).$$

Recall that $A^{[j_k]}$ is the j_k 'th column of A , while $B_{[j_k]}$ is the j_k 'th row of B ; in particular, every determinant appearing in the formula is the determinant of an $m \times m$ matrix. If

$m = n$, then J consists of the single list $[1, 2, \dots, n]$, and the Cauchy–Binet formula reduces to the product formula for square matrices.

To prove the Cauchy–Binet formula, we imitate the proof of the product formula from the last section. First, the definitions of determinants and matrix products give

$$\det(AB) = \sum_{f \in S_m} \operatorname{sgn}(f) \prod_{i \in [m]} (AB)(f(i), i) = \sum_{f \in S_m} \operatorname{sgn}(f) \prod_{i \in [m]} \sum_{k_i \in [n]} A(f(i), k_i) B(k_i, i).$$

Using the distributive law repeatedly (see Exercise 5), we get

$$\det(AB) = \sum_{f \in S_m} \sum_{k_1 \in [n]} \cdots \sum_{k_m \in [n]} \operatorname{sgn}(f) \prod_{i \in [m]} A(f(i), k_i) \prod_{i \in [m]} B(k_i, i).$$

Now we move the summation over f inside the other sums, replace the multiple summations over indices k_1, \dots, k_m by a single summation over lists $k = [k_1, \dots, k_m] \in [n]^m$, and rearrange factors in the resulting summation:

$$\det(AB) = \sum_{k \in [n]^m} \prod_{i \in [m]} B(k_i, i) \left[\sum_{f \in S_m} \operatorname{sgn}(f) \prod_{i \in [m]} A(f(i), k_i) \right].$$

We recognize the term in brackets as the definition of $\det(A^{[k_1]}, \dots, A^{[k_m]})$. If $k_i = k_j$ for some $i \neq j$, the matrix in question has two equal columns, so its determinant is zero. Discarding these terms, we now have

$$\det(AB) = \sum_{\substack{k \in [n]^m \\ k_i \text{ distinct}}} \det(A^{[k_1]}, \dots, A^{[k_m]}) \prod_{i \in [m]} B(k_i, i).$$

Given a list $k = [k_1, \dots, k_m]$ of distinct elements of $[n]$, let $\operatorname{sort}(k)$ be the list obtained by sorting k into increasing order. Note that $j = \operatorname{sort}(k)$ lies in J . Defining $\operatorname{inv}(k)$ and $\operatorname{sgn}(k)$ as in the case $m = n$ (see §2.6), we see that it takes $\operatorname{inv}(k)$ interchanges of adjacent columns to turn the matrix $(A^{[k_1]}, \dots, A^{[k_m]})$ into the matrix $(A^{[j_1]}, \dots, A^{[j_m]})$. Therefore, $\det(A^{[k_1]}, \dots, A^{[k_m]}) = \operatorname{sgn}(k) \det(A^{[j_1]}, \dots, A^{[j_m]})$. Grouping summands in the formula for $\det(AB)$ based on the value of $j = \operatorname{sort}(k)$, we obtain

$$\begin{aligned} \det(AB) &= \sum_{j \in J} \sum_{\substack{k \in [n]^m \\ \operatorname{sort}(k)=j}} \det(A^{[k_1]}, \dots, A^{[k_m]}) \prod_{i \in [m]} B(k_i, i) \\ &= \sum_{j \in J} \det(A^{[j_1]}, \dots, A^{[j_m]}) \left[\sum_{\substack{k \in [n]^m \\ \operatorname{sort}(k)=j}} \operatorname{sgn}(k) \prod_{i \in [m]} B(k_i, i) \right]. \end{aligned}$$

Finally, compare the term in square brackets to the definition of $\det(B_{[j_1]}, \dots, B_{[j_m]})$. The formula in brackets looks just like the original definition of the determinant of a matrix, except now we are indexing the rows of the matrix by the increasing sequence $j_1 < \dots < j_m$ instead of the standard indexing sequence $1 < 2 < \dots < m$. This renaming of indices makes no difference when we calculate $\operatorname{inv}(k)$ or $\operatorname{sgn}(k)$. Therefore, the term in brackets is none other than $\det(B_{[j_1]}, \dots, B_{[j_m]})$. (To prove this formally, one must change the summation variable from a sum over the set of bijections $k : [m] \rightarrow \{j_1, \dots, j_m\}$ to a sum over the set S_m of bijections from $[m]$ to $[m]$.) In summary,

$$\det(AB) = \sum_{j \in J} \det(A^{[j_1]}, A^{[j_2]}, \dots, A^{[j_m]}) \det(B_{[j_1]}, B_{[j_2]}, \dots, B_{[j_m]}).$$

For another proof, see Exercise 49 in Chapter 20.

5.15 Cayley–Hamilton Theorem

Given $A \in M_n(R)$, the *characteristic polynomial* of A , denoted χ_A , is defined by $\chi_A = \det(xI_n - A) \in R[x]$. This section presents a rather tricky proof of the famous *Cayley–Hamilton theorem*, which states that $\chi_A(A) = 0$ for all $A \in M_n(R)$. In more detail, if $\chi_A = \sum_{i \geq 0} c_i x^i \in R[x]$ for some $c_i \in R$, then the theorem asserts that $\sum_{i \geq 0} c_i A^i$ is the zero matrix. Other, potentially more intuitive proofs of the Cayley–Hamilton theorem will be given in §8.14 and §18.19.

Before giving the official proof, we need some preliminary discussion. For fixed $A \in M_n(R)$, an *eigenvalue* of A is a scalar $c \in R$ such that there exists a nonzero $x \in R^n$ with $Ax = cx$; any such x is called an *eigenvector* associated with c . Let $p = \det(xI_n - A)$ be the characteristic polynomial of A . Informally, since p is a polynomial in the “variable” x with coefficients in the commutative ring R , we can evaluate this polynomial at $x = \lambda$, for any $\lambda \in R$. Doing this produces $p(\lambda) = \det(\lambda I_n - A)$. Suppose for a moment that R is a field. In this case, $\lambda \in R$ is a root of p iff $p(\lambda) = 0$ iff $\det(\lambda I_n - A) = 0$ iff $\lambda I_n - A$ is a non-invertible matrix iff (by §4.13) there exists a nonzero $x \in R^n$ with $(\lambda I_n - A)x = 0$ iff there exists a nonzero $x \in R^n$ with $Ax = \lambda x$ iff λ is an eigenvalue of A . So, we obtain the result that *the eigenvalues of a square matrix with entries in a field are precisely the roots of the characteristic polynomial of that matrix*.

Returning to the case of a commutative ring R , suppose we “evaluate the polynomial $p = \det(xI_n - A)$ at $x = A$.” We seem to obtain $p(A) = \det(AI_n - A) = \det(0_n) = 0$, giving the conclusion of the Cayley–Hamilton theorem. However, something must be wrong with this calculation, since $p(A)$ is an $n \times n$ matrix, but $\det(0_n) = 0_R$ is a scalar! So, we need to be more careful about “evaluating” the variable x appearing in the characteristic polynomial. To do so, we must recall the *universal mapping property (UMP) for polynomials* proved in §3.5: given two commutative rings R and S , a ring homomorphism $h : R \rightarrow S$, and $c \in S$, there exists a unique ring homomorphism $H : R[x] \rightarrow S$ such that $H(r) = h(r)$ for all $r \in R$, and $H(x) = c$. In our discussion of eigenvalues in the previous paragraph, we took $S = R$, $c = \lambda$, and h to be the identity map on R ; then $H(p) = p(\lambda)$ for all $p \in R[x]$.

In the current setting, where we want to evaluate x at A , we cannot use the same h since A does not belong to R . In fact, A belongs to the *non-commutative* ring $M_n(R)$. To get around this, we define S to be the set of all matrices in $M_n(R)$ of the form $r_0 I_n + r_1 A + r_2 A^2 + \cdots + r_k A^k$ for some $k \geq 0$ and $r_0, r_1, \dots, r_k \in R$. The reader should check that S is a subring of $M_n(R)$, and that S is *commutative*. We define $h : R \rightarrow S$ by setting $h(r) = rI_n$ for $r \in R$; one readily checks that h is a ring homomorphism. The UMP tells us there is a ring homomorphism $H : R[x] \rightarrow S$ extending h and sending x to A . More specifically, for any $r_0, \dots, r_k \in R$,

$$H(r_0 + r_1 x + \cdots + r_k x^k) = r_0 I_n + r_1 A + \cdots + r_k A^k \in S.$$

Now, finally, we have a precise way of talking about $p(A)$: namely, $p(A) = H(p)$ is the image of p under the ring homomorphism H that extends h and sends x to A .

Since $p = \det(xI_n - A)$, we need to find $H(p) = H(\det(xI_n - A))$. For this, we need another technical digression. Suppose $f : T \rightarrow U$ is a ring homomorphism between commutative rings T and U . We claim that there is an associated ring homomorphism $F : M_n(T) \rightarrow M_n(U)$ such that $F(B)_{i,j} = f(B_{i,j})$ for all $B \in M_n(T)$ and all $i, j \in [n]$. This definition says that F acts on a matrix B by applying f to each entry $B(i, j)$. We claim also that $f(\det(B)) = \det(F(B))$ for all $B \in M_n(T)$. The reader is asked to prove these two claims in Exercise 3.

We apply these claims to the ring homomorphism $H : R[x] \rightarrow S$, obtaining a ring

homomorphism $H^* : M_n(R[x]) \rightarrow M_n(S)$ that acts by applying H to every entry of a matrix in $M_n(R[x])$. Let $C = H^*(xI_n - A)$; then $\det(C) = \det(H^*(xI_n - A)) = H(\det(xI_n - A)) = H(p)$. Note that $C \in M_n(S)$ is an $n \times n$ matrix, each entry of which is itself an $n \times n$ matrix built up from powers of A .

For example, let $n = 2$ and $A = \begin{bmatrix} r & s \\ t & u \end{bmatrix}$ where $r, s, t, u \in R$. Then $B = \begin{bmatrix} x-r & -s \\ -t & x-u \end{bmatrix}$ and $p = \chi_A = \det(B) = (x-r)(x-u) - st = x^2 - (r+u)x + (ru-st)$, so

$$p(A) = A^2 - (r+u)A + (ru-st)I_2.$$

On the other hand, applying H to every entry of B produces the matrix of matrices

$$C = \begin{bmatrix} A - rI_2 & -sI_2 \\ -tI_2 & A - uI_2 \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} 0 & s \\ t & u-r \end{bmatrix} & \begin{bmatrix} -s & 0 \\ 0 & -s \end{bmatrix} \\ \begin{bmatrix} -t & 0 \\ 0 & -t \end{bmatrix} & \begin{bmatrix} r-u & s \\ t & 0 \end{bmatrix} \end{bmatrix}.$$

We see $\det(C) = (A - rI_2)(A - uI_2) - (-sI_2)(-tI_2) = A^2 - (r+u)A + (ru-st)I_2 = p(A)$. For general n and A , the definitions show that, for all $i, j \in [n]$ with $j \neq i$, $C(i, i) = A - A(i, i)I_n$ and $C(i, j) = -A(i, j)I_n$. We know that $p(A) = H(p) = \det(C)$, so our task is to prove that $\det(C)$ is the zero matrix in S . For reasons that will emerge shortly, we will actually prove the equivalent statement $\det(C^T) = 0$.

For an ordinary matrix $A \in M_n(R)$, we can compute the vector-matrix product wA for any row vector $w \in R^n$ via the formula $(wA)_j = \sum_{i=1}^n w_i A(i, j)$ for $j \in [n]$. Since $C^T \in M_n(S)$ is a matrix of matrices, we will instead compute an expression of the form zC^T , where z is a row vector of row vectors. More precisely, suppose $U \in M_n(S)$ and $z = [z_1 \ z_2 \ \cdots \ z_n] \in (R^n)^n$, where each $z_i \in R^n$ is a row vector. Then $zU \in (R^n)^n$ is the row vector whose j 'th component is the row vector $\sum_{i=1}^n z_i U(i, j)$. We let the reader verify the *associativity property* $(zU)V = z(UV)$, valid for all $z \in (R^n)^n$ and $U, V \in M_n(S)$.

Continuing the proof, let us take $z = [e_1 \ e_2 \ \cdots \ e_n]$, where each e_i is the row vector of length n with a 1 in position i and zeroes elsewhere. We claim $zC^T = 0$, where 0 denotes the row vector consisting of n row vectors all equal to zero. To see this, note the j 'th row vector appearing in zC^T is

$$\begin{aligned} \sum_{i=1}^n e_i C^T(i, j) &= e_j C(j, j) + \sum_{i \neq j} e_i C(j, i) \\ &= e_j (A - A(j, j)I_n) + \sum_{i \neq j} e_i (-A(j, i)I_n) = e_j A - \sum_{i=1}^n A(j, i)e_i. \end{aligned}$$

Note that $e_j A$ and $\sum_{i=1}^n A(j, i)e_i$ are both equal to the j 'th row of A , so their difference is indeed zero. This holds for all j , so $zC^T = 0$ as claimed.

To finish, multiply the equation $zC^T = 0$ on the right by the classical adjoint $\text{adj}(C^T)$. The right side is still $0 \in (R^n)^n$, and the left side becomes $(zC^T)\text{adj}(C^T) = z(C^T \text{adj}(C^T)) = z(\det(C^T)I_n)$. Using the definition of zU one more time, we see that the n row vectors appearing in $0 = z(\det(C^T)I_n)$ are $e_1 \det(C^T), \dots, e_n \det(C^T)$, which are precisely the n rows of $\det(C^T)$. Since all these row vectors are zero, the matrix $\det(C^T) = \det(C)$ must be zero, completing the proof of the Cayley–Hamilton theorem.

Continuing the previous example where $n = 2$, we have

$$zC^T = [[1 \ 0] [0 \ 1]] \begin{bmatrix} A - rI_2 & -tI_2 \\ -sI_2 & A - uI_2 \end{bmatrix}.$$

The first component of the answer is $[1 \ 0](A - rI_2) + [0 \ 1](-sI_2) = [r \ s] - [r \ 0] - [0 \ s] = [0 \ 0]$. Similarly, the second component of the answer is $[1 \ 0](-tI_2) + [0 \ 1](A - uI_2) = [-t \ 0] + [t \ u] - [0 \ u] = [0 \ 0]$. Multiplying on the right by $\text{adj}(C^T)$ gives

$$[[0 \ 0] \ [0 \ 0]] = (zC^T)\text{adj}(C^T) = z(\det(C^T)I_2) = [[1 \ 0] \ [0 \ 1]] \begin{bmatrix} \det(C) & 0 \\ 0 & \det(C) \end{bmatrix}.$$

Multiplying out the far right side shows that row 1 and row 2 of $\det(C)$ are $[0 \ 0]$, so $\det(C) = 0$.

5.16 Permanents

The *permanent* of a square matrix $A \in M_n(R)$ is defined by the formula

$$\text{per}(A) = \sum_{f \in S_n} \prod_{j \in [n]} A(f(j), j).$$

This formula is just like the defining formula (5.1) for $\det(A)$, except that the sign $\text{sgn}(f)$ has been omitted. Note that the permanent of A is the sum of $n!$ terms, each of which is a product of n elements of A chosen from distinct rows and columns in all possible ways.

Many of the results in this chapter extend without difficulty to permanents; one merely deletes all occurrences of $\text{sgn}(f)$ in the relevant proofs. For example, the changes of variable $g = f^{-1}$ followed by $j = g(i)$ show that

$$\text{per}(A) = \sum_{g \in S_n} \prod_{i \in [n]} A(i, g(i)) = \text{per}(A^T).$$

Similarly, one establishes the following facts about permanents: $\text{per}(\bar{A}) = \text{per}(A^*) = \overline{\text{per}(A)}$ when $R = \mathbb{C}$; $\text{per} : (R^n)^n \rightarrow R$ is an R -multilinear function of the n rows of A (or the n columns of A); the permanent of a matrix with a row (or column) of zeroes is zero; the permanent of a triangular or diagonal matrix is the product of the diagonal entries; multiplying one row (or column) of A by $c \in R$ multiplies the permanent by c ; multiplying the whole matrix by c multiplies the permanent by c^n .

Instead of the alternating property of the determinant, the permanent has the following *symmetry property*: if the rows (or columns) of A are permuted in any fashion, the permanent is unchanged. To prove this, it suffices to show $\text{per}(B) = \text{per}(A)$ where B is obtained from A by interchanging row i and row j , since any permutation of rows can be achieved by a sequence of such interchanges. Using the change of variable $g = f \circ (i, j)$ (so $f = g \circ (i, j)$), we find that

$$\begin{aligned} \text{per}(A) &= \sum_{f \in S_n} \prod_{k \in [n]} A(k, f(k)) = \sum_{f \in S_n} A(i, f(i))A(j, f(j)) \prod_{k \neq i, j} A(k, f(k)) \\ &= \sum_{g \in S_n} A(i, g(j))A(j, g(i)) \prod_{k \neq i, j} A(k, g(k)) = \sum_{g \in S_n} B(j, g(j))B(i, g(i)) \prod_{k \neq i, j} B(k, g(k)) \\ &= \sum_{g \in S_n} \prod_{k \in [n]} B(k, g(k)) = \text{per}(B). \end{aligned}$$

Note that it is *no longer true* that the permanent of a matrix with two equal rows

(or columns) is always zero. For example, the permanent of $\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ is $1 \cdot 1 + 1 \cdot 1 = 2$. In other words, the permanent function is *not alternating*. Consequently, those properties of determinants whose proofs invoked the alternating property are no longer valid for permanents. For instance, if B is obtained from A by adding c times row j of A to row i of A , we *cannot* conclude that $\text{per}(B) = \text{per}(A)$. Unfortunately, this means that we can no longer use Gaussian elimination as an efficient method for computing permanents as we did for determinants. The Laplace expansions for permanents still hold if we omit the signs; for instance, for any $i \in [n]$ we have

$$\text{per}(A) = \sum_{k=1}^n A(i, k) \text{per}(A[i|k]).$$

(The general case can be deduced from the special case $i = n$ by using the symmetry of the permanent to reorder rows.) However, the recursive Laplace expansion still requires the evaluation of $n!$ terms to compute the permanent of A . The formula $AA' = \det(A)I_n = A'A$ does not generalize to permanents, nor does the product formula $\text{per}(AB) = \text{per}(A)\text{per}(B)$ hold in general.

5.17 Summary

Here are the main facts we proved about the determinant of a square $n \times n$ matrix A with entries in a commutative ring R .

1. *Two Formulas for the Determinant:*

$$\det(A) = \sum_{f \in S_n} \text{sgn}(f) \prod_{i=1}^n A(f(i), i) = \sum_{f \in S_n} \text{sgn}(f) \prod_{i=1}^n A(i, f(i)).$$

2. *Transpose and Conjugation:* $\det(A^T) = \det(A)$, $\det(\bar{A}) = \overline{\det(A)} = \det(A^*)$.

3. *Diagonal and Triangular Matrices:* If A is upper-triangular, lower-triangular, or diagonal, then $\det(A) = \prod_{i=1}^n A(i, i)$. In particular, $\det(I_n) = 1_R$ and $\det(0_n) = 0_R$.

4. *Multilinearity and the Alternating Property:* We can view $\det : (R^n)^n \rightarrow R$ as a function of the n rows of A (or the n columns of A): $\det(A) = \det(A_{[1]}, \dots, A_{[n]}) = \det(A^{[1]}, \dots, A^{[n]})$. Either way, \det is an R -multilinear function of its n arguments. In other words, if all rows other than row i are fixed, then the function $v \in R^n \mapsto \det(A_1, \dots, A_{i-1}, v, A_{i+1}, \dots, A_n)$ is R -linear; similarly for columns. If a matrix A has two equal rows, two equal columns, a zero row, or a zero column, then $\det(A) = 0$.

5. *Elementary Operations and Determinants:* There are three elementary row and column operations on matrices, which may be effected by multiplying on the left or right by appropriate elementary matrices. First, multiplying one row of a matrix by $c \in R$ multiplies the determinant by c ; the corresponding elementary matrix has determinant c (this matrix is elementary only for c invertible). Second, interchanging any two distinct rows of a matrix multiplies the determinant by -1_R ; the corresponding elementary matrix has determinant -1_R . Third, adding a scalar multiple of one row of a matrix to another row leaves the determinant unchanged; the corresponding elementary matrix has determinant 1_R . Combining these facts with row-reduction via Gaussian elimination, we obtain an efficient method for computing determinants when R is a field.

6. *Product Formula:* We have $\det(AB) = \det(A)\det(B)$ for all $A, B \in M_n(R)$.
 7. *Laplace Expansions:* For all $i, j \in [n]$, we have

$$\det(A) = \sum_{j=1}^n A(i, j)(-1)^{i+j} \det(A[i|j]) = \sum_{i=1}^n A(i, j)(-1)^{i+j} \det(A[i|j]).$$

Here $A[i|j]$ is the matrix in $M_{n-1}(R)$ obtained by deleting row i and column j of A .

8. *Classical Adjoint:* The classical adjoint of $A \in M_n(R)$ is the matrix $A' = \text{adj}(A)$ whose i, j -entry is $(-1)^{i+j} \det(A[j|i])$. We have $A'A = \det(A)I_n = AA'$.
 9. *Explicit Formula for Inverses:* A matrix $A \in M_n(R)$ is invertible iff $\det(A)$ is an invertible element in the commutative ring R . In this case, $A^{-1} = \det(A)^{-1} \text{adj}(A)$, so that

$$A^{-1}(i, j) = \det(A)^{-1}(-1)^{i+j} \det(A[j|i]).$$

10. *Cramer's Rule:* If $A\vec{x} = \vec{b}$ is a system of n linear equations in n unknowns whose coefficient matrix A is invertible, then the unique solution is given by $x_i = \det(A)^{-1} \det(C_i)$ for each $i \in [n]$, where C_i is the matrix obtained from A by substituting \vec{b} for the i 'th column of A .

11. *Cauchy–Binet Formula:* If $A \in M_{m,n}(R)$, $B \in M_{n,m}(R)$, and $m \leq n$, then

$$\det(AB) = \sum_{1 \leq j_1 < j_2 < \dots < j_m \leq n} \det(A^{[j_1]}, A^{[j_2]}, \dots, A^{[j_m]}) \det(B_{[j_1]}, B_{[j_2]}, \dots, B_{[j_m]}).$$

12. *Characteristic Polynomials:* The characteristic polynomial of $A \in M_n(R)$ is $\chi_A = \det(xI_n - A)$. Writing $\chi_A = \sum_{i \geq 0} c_i x^i$ with $c_i \in R$, the Cayley–Hamilton theorem states that $\chi_A(A) = \sum_{i \geq 0} c_i A^i = 0$. When R is a field, the roots of χ_A in R are precisely the eigenvalues of A .

13. *Permanents:* For $A \in M_n(R)$, we have

$$\text{per}(A) = \sum_{f \in S_n} \prod_{j \in [n]} A(f(j), j) = \sum_{f \in S_n} \prod_{i \in [n]} A(i, f(i)) = \text{per}(A^T).$$

The permanent can be viewed as a symmetric, R -multilinear function $\text{per} : (R^n)^n \rightarrow R$ whose n arguments are the rows (or columns) of A . We have $\text{per}(A) = \prod_{i \in [n]} A(i, i)$ when A is triangular or diagonal, and the Laplace expansions

$$\text{per}(A) = \sum_{j=1}^n A(i, j) \text{per}(A[i|j]) = \sum_{i=1}^n A(i, j) \text{per}(A[i|j])$$

are valid. However, the permanent is not alternating; we cannot use elementary row operations to evaluate permanents; and it is not true that $\text{per}(AB) = \text{per}(A)\text{per}(B)$ in general.

5.18 Exercises

Unless otherwise specified, assume R is a commutative ring in these exercises.

- (a) Use (5.1) to write a formula for $\det(A)$ when $A \in M_2(R)$. (b) Repeat (a) for $A \in M_4(R)$.
- Suppose $A \in M_5(R)$ satisfies $A(i, j) = 0_R$ for all $i, j \in [5]$ with $i + j$ odd. Use (5.1) to compute $\det(A)$.
- Suppose R and S are commutative rings, and $f : R \rightarrow S$ is a ring homomorphism.
 (a) Prove: for all $n \geq 1$, there is a ring homomorphism $F : M_n(R) \rightarrow M_n(S)$ given by $F(A)_{i,j} = f(A_{i,j})$ for all $A \in M_n(R)$ and all $i, j \in [n]$. (b) Prove: for all $A \in M_n(R)$, $\det(F(A)) = f(\det(A))$. (c) Point out why the formula $\det(\overline{A}) = \overline{\det(A)}$ (for $A \in M_n(\mathbb{C})$) is a special case of (b).
- Suppose $A \in M_n(R[x])$ has the property that for all $i, j \in [n]$, $A(i, j)$ is either zero or has degree at most d_i . Prove $\det(A) \in R[x]$ is either zero or has degree at most $d_1 + d_2 + \dots + d_n$.
- Let R be a ring, and let $a_{ij} \in R$ for $1 \leq i \leq m$ and $1 \leq j \leq n_i$. Use the distributive axioms for R and induction to prove the *generalized distributive law*

$$\prod_{i=1}^m \left(\sum_{j_i=1}^{n_i} a_{ij_i} \right) = \sum_{j_1=1}^{n_1} \sum_{j_2=1}^{n_2} \dots \sum_{j_m=1}^{n_m} \left(\prod_{i=1}^m a_{ij_i} \right).$$

(This formula was used several times in §5.13 and §5.14.)

- Let $A \in M_n(R)$. (a) Use (5.1) to show that the characteristic polynomial χ_A is a monic polynomial of degree n . (b) Prove that the coefficient of x^{n-1} in $\chi(A)$ is $-\text{tr}(A) = -\sum_{i=1}^n A(i, i)$. (c) Prove that the constant coefficient in $\chi(A)$ is $(-1)^n \det(A)$.

- (a) Compute $\det \begin{bmatrix} 0 & 4 & 0 & 0 & 0 \\ 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 5 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \end{bmatrix}$. (b) Describe how to find the determinant of a matrix that has exactly one nonzero entry in every row and column.

- (a) Prove in detail that for an upper-triangular $A \in M_n(R)$, $\det(A) = \prod_{i=1}^n A(i, i)$. (b) Suppose $A \in M_n(R)$ satisfies $A(i, j) = 0_R$ for all $i, j \in [n]$ with $i + j > n + 1$. Find and prove a simple formula for $\det(A)$.
- Suppose $A \in M_n(R)$ satisfies $A(i, j) = 0$ for all $i, j \in [n]$ with $i > j + 1$. In general, how many of the $n!$ terms in (5.1) might be nonzero for such an A ?
- (a) Consider a matrix $U \in M_n(R)$ written in “block form” $U = \begin{bmatrix} A & B \\ 0 & D \end{bmatrix}$, where $A \in M_k(R)$, $B \in M_{k, n-k}(R)$, $0 \in M_{n-k, k}(R)$, and $D \in M_{n-k}(R)$. Prove carefully that $\det(U) = \det(A)\det(D)$. (b) Suppose $k = n - k$ and we replace the zero block in (a) by $C \in M_k(R)$. Prove or disprove: $\det(U) = \det(A)\det(D) - \det(C)\det(B)$.
- Call $U \in M_n(R)$ *block upper-triangular* iff there exist positive integers s, k_1, \dots, k_s such that $k_1 + \dots + k_s = n$ and for all integers i, j, t with $k_1 + \dots + k_{t-1} < i \leq k_1 + \dots + k_t$ and all $j \leq k_1 + \dots + k_{t-1}$, $U(i, j) = 0$. Draw a picture of such a

- matrix, and derive a formula for its determinant. (Use Exercise 10 and induction on s .)
12. Let (G, \star) be any group. For $g \in G$, define maps L_g , R_g , and C_g from G to G by setting $L_g(x) = g \star x$, $R_g(x) = x \star g$, and $C_g(x) = g \star x \star g^{-1}$ for all $x \in G$. (a) Prove L_g is a bijection by showing L_g is one-to-one and onto. (b) Prove R_g is a bijection by showing $R_{g^{-1}}$ is a two-sided inverse of R_g . (c) Prove C_g is a group isomorphism. (d) Prove $I : G \rightarrow G$, given by $I(g) = g^{-1}$ for $g \in G$, is a bijection. (e) If T is any ring, G is finite, and $f : G \rightarrow T$ is any function, use L_g to explain why $\sum_{x \in G} f(x) = \sum_{x \in G} f(g \star x)$ for each fixed $g \in G$. Write similar formulas based on the maps R_g , C_g , and I .
 13. (a) Assume $A \in M_n(\mathbb{R})$ satisfies $AA^T = I_n$. Prove $\det(A) \in \{+1, -1\}$. (b) Must (a) be true if \mathbb{R} is replaced by an arbitrary commutative ring R ? Prove or give a counterexample. (c) Assume $A \in M_n(\mathbb{C})$ satisfies $AA^* = I_n$. Prove $|\det(A)| = 1$ and $A^*A = I_n$.
 14. (a) Assume n is odd and $A \in M_n(\mathbb{R})$ satisfies $A^T = -A$. Prove $\det(A) = 0$. (b) Must (a) be true if \mathbb{R} is replaced by an arbitrary commutative ring R ? Prove or give a counterexample.
 15. (a) Assume $A \in M_n(\mathbb{R})$ satisfies $A^k = 0$ for some $k > 0$. Prove $\det(A) = 0$. (b) Must (a) be true if \mathbb{R} is replaced by an arbitrary commutative ring R ? Prove or give a counterexample.
 16. Use Laplace expansions and induction on n to reprove that $\det(A) = \det(A^T)$ for all $A \in M_n(R)$.
 17. (a) Fix row vectors $A_1 = (2, 6, -1, -1)$, $A_2 = (3, 4, 0, -2)$, and $A_4 = (2, -3, 1, 5)$ in \mathbb{R}^4 . Define $D : \mathbb{R}^4 \rightarrow \mathbb{R}$ by $D(v) = \det(A_1, A_2, v, A_4)$ for $v = (v_1, v_2, v_3, v_4) \in \mathbb{R}^4$. Write $D(v)$ explicitly in the form $c_1v_1 + \dots + c_4v_4$ for some $c_i \in \mathbb{R}$, and confirm that this is an \mathbb{R} -linear function. (b) Repeat (a) using column vectors $A^{[2]} = [3 \ 5 \ 4]^T$, $A^{[3]} = [2 \ 7 \ -1]^T$, and $D : \mathbb{R}^3 \rightarrow R$ defined on column vectors $v = [v_1 \ v_2 \ v_3]^T \in \mathbb{R}^3$ by $D(v) = \det(v, A^{[2]}, A^{[3]})$.
 18. In the formula $D(v) = \sum_{j=1}^n c_j v(j)$ from §5.6, find an expression for c_j involving a sign times a determinant of a submatrix of A .
 19. (a) Use elementary row operations to prove that if some row of $A \in M_n(R)$ is an R -linear combination of other rows of A , then $\det(A) = 0_R$. (b) Use (a) to prove that if some column of $A \in M_n(R)$ is an R -linear combination of other columns of A , then $\det(A) = 0_R$. (c) Suppose the rows of $A \in M_n(R)$ are R -linearly dependent, i.e., there exist $r_1, \dots, r_n \in R$ (not all zero) with $r_1 A_{[1]} + \dots + r_n A_{[n]} = 0$ in R^n . Give a specific example to show that $\det(A)$ need not be zero.
 20. (a) Use Laplace expansions to reprove the fact that multiplying row i of $A \in M_n(R)$ by $c \in R$ replaces $\det(A)$ by $c \det(A)$. (b) Use Laplace expansions to reprove the fact that interchanging two rows of a matrix changes the sign of the determinant. (c) Use Laplace expansions and induction to reprove the fact that adding c times row i to row $j \neq i$ of $A \in M_n(R)$ does not change $\det(A)$.
 21. Let A_1, \dots, A_n be row vectors in R^n , and let $c_{ij} \in R$ for each $i < j$ in $[n]$. Prove that $\det(A_1, \dots, A_n)$ equals

$$\det \left(A_1 + \sum_{j>1} c_{1j} A_j, A_2 + \sum_{j>2} c_{2j} A_j, \dots, A_i + \sum_{j>i} c_{ij} A_j, \dots, A_n \right)$$

- (a) by using elementary row operations; (b) as a special case of the formula $\det(BA) = \det(B)\det(A)$.
22. (a) Prove: for all elementary matrices $E \in M_n(R)$, E^T is also an elementary matrix. (b) Prove or disprove: for all elementary matrices $E \in M_n(R)$ and all $k \in \mathbb{Z}$, E^k is also an elementary matrix.
23. (a) Use the definition (5.1) to prove that $\det(A)$ is a multilinear function of the n columns of $A \in M_n(R)$. (b) Prove that each elementary column operation has the expected effect on determinants, without invoking the corresponding facts about elementary row operations. (c) Use (b) to prove that $\det(AE) = \det(A)\det(E)$ for all $A \in M_n(R)$ and all elementary matrices $E \in M_n(R)$.
24. Given $a, b \in R$, evaluate $\det(aI_n + bJ_n)$, where $J_n \in M_n(R)$ has every entry equal to 1_R .
25. Let $n \in \mathbb{N}^+$ and define $A \in M_n(\mathbb{R})$ by $A(i, j) = \min(i, j)$ for $i, j \in [n]$. Compute $\det(A)$.
26. (a) Given distinct r_1, r_2, \dots, r_n in a field F , define a matrix $A \in M_{n+1}(F[x])$ by setting $A(i, i) = x$ for $i \in [n]$; $A(i, n+1) = 1$ for $i \in [n+1]$; $A(i, j) = r_j$ for $i > j$ in $[n+1]$; and $A(i, j) = r_{j-1}$ for $i < j$ in $[n]$. Prove $\det(A) = (x - r_1)(x - r_2) \cdots (x - r_n)$ by showing that the determinant is a monic polynomial of degree n having every r_j as a root. (You may need to apply Exercise 3 to certain evaluation homomorphisms.) (b) Suppose r_1, r_2, \dots, r_n are not necessarily distinct elements of a commutative ring R . Does the formula for $\det(A)$ in (a) still hold?
27. Let $A \in M_n(F)$ satisfy $A(i, j) = 0$ for all $i \neq j$ in $[n-1]$. Compute $\det(A)$.
28. Compute the determinant of each matrix by using row operations to reduce the matrix to an upper-triangular matrix.

$$(a) A = \begin{bmatrix} 1 & 4 & 6 \\ 3 & 2 & 1 \\ -2 & 4 & 5 \end{bmatrix} \in M_3(\mathbb{R}); \quad (b) B = \begin{bmatrix} 4 & 0 & 1 & 2 \\ 3 & 3 & 1 & 4 \\ 2 & 1 & 1 & 3 \\ 3 & 1 & 3 & 4 \end{bmatrix} \in M_4(\mathbb{Z}_5);$$

$$(c) C = \begin{bmatrix} 1 & 1 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 1 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \in M_7(\mathbb{Z}_2).$$

29. Compute the determinants of each matrix in Exercise 28 by repeated use of Laplace expansions along convenient rows and columns.
30. *Vandermonde Determinant.* Given $n \in \mathbb{N}^+$ and $x_1, \dots, x_n \in R$, define $V \in M_n(R)$ by $V(i, j) = x_i^{n-j}$ for $i, j \in [n]$. Prove $\det(V) = \prod_{1 \leq i < j \leq n} (x_i - x_j)$. (Use elementary row and column operations and induction on n .)
31. (a) Given $n \geq 0$, distinct elements a_0, \dots, a_n in a field F , and $b_0, \dots, b_n \in F$, prove there exist unique $c_0, \dots, c_n \in F$ such that the polynomial $p = \sum_{i=0}^n c_i x^i \in F[x]$ satisfies $p(a_i) = b_i$ for $0 \leq i \leq n$. Do this by setting up a system of $n+1$ linear equations in $n+1$ unknowns, and using the Vandermonde matrix V and its determinant from Exercise 30. (b) Contemplate the relationship between (a) and

the Lagrange interpolation formula (see §3.18). Can you use the latter formula to describe the columns of V^{-1} ?

32. Given $p = a_0 + a_1x + a_2x^2 + \cdots + a_{n-1}x^{n-1} + x^n \in R[x]$, define the *companion matrix* $C_p \in M_n(R)$ by setting $C_p(j+1, j) = 1$ for $j \in [n-1]$, $C_p(i, n) = -a_{i-1}$ for $i \in [n]$, and letting all other entries be zero. Prove that $p = \det(xI_n - C_p) = \chi_{C_p}$.
33. Given $S \subseteq [n]$, let $\text{inv}(S)$ be the number of ordered pairs (s, t) with $s \in S$, $t \in [n] \sim S$, and $s > t$, and let $\text{sgn}(S) = (-1)^{\text{inv}(S)}$. (a) Fix $n \in \mathbb{N}^+$ and a k -element subset I of $[n]$. Prove the *generalized Laplace expansion*: for all $A \in M_n(R)$,

$$\det(A) = \text{sgn}(I) \sum_{\substack{J \subseteq [n] \\ |J|=k}} \text{sgn}(J) \det(A[I|J]) \det(A[[n] \sim I|[n] \sim J]).$$

In this formula, $A[I|J]$ denotes the matrix in $M_{n-k}(R)$ obtained by erasing all rows in I and all columns in J , and $A[[n] \sim I|[n] \sim J]$ denotes the matrix in $M_k(R)$ obtained by erasing all rows not in I and all columns not in J . (Reduce to the case $I = [k]$.) (b) In the formula in (a), show that $\text{sgn}(I) \text{sgn}(J) = (-1)^{\sum_{i \in I} i + \sum_{j \in J} j}$. (c) Show how to deduce the Laplace expansion of $\det(A)$ along row i as a special case of (a).

34. Use Laplace expansions along columns to prove that $A'A = (\det A)I_n$ for all $A \in M_n(R)$.
35. (a) Given $a, b, c, d \in R$, find the classical adjoint of $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$. (b) Use (a) to write an explicit formula for A^{-1} , when it exists.
36. (a) Find the classical adjoint and (if possible) the inverse of

$$A = \begin{bmatrix} 4 & 1 & -2 \\ 3 & 0 & 5 \\ 7 & -1 & 3 \end{bmatrix} \in M_3(\mathbb{R}).$$

(b) Find the classical adjoint and (if possible) the inverse of

$$B = \begin{bmatrix} 2 & 3 & 1 & 0 \\ 1 & 5 & 3 & 1 \\ 2 & 2 & 6 & 4 \\ 3 & 4 & 5 & 1 \end{bmatrix} \in M_4(\mathbb{Z}_7).$$

37. Use the explicit formula for A^{-1} to prove that if A is a unitriangular matrix (see Exercise 20 of Chapter 4), then A^{-1} exists and is also unitriangular.

38. Use Cramer's rule to find the solutions to $A\vec{x} = \vec{b}$, where $A = \begin{bmatrix} 4 & 1 & -2 \\ 3 & 6 & -1 \\ 1 & 1 & 5 \end{bmatrix}$

and $\vec{b} = [b_1 \ b_2 \ b_3]^T$ for some $b_1, b_2, b_3 \in \mathbb{R}$.

39. Consider the linear system of equations

$$\left\{ \begin{array}{rclclclclcl} 2x_1 & +3x_2 & -x_3 & & -x_5 & = & 4 \\ 3x_1 & +5x_2 & +x_3 & +2x_4 & +x_5 & = & 0 \\ -x_1 & & & -3x_4 & +x_5 & = & 1 \\ & x_2 & +4x_3 & +x_4 & & = & -1 \\ & -2x_2 & -3x_3 & & +4x_5 & = & 6 \end{array} \right.$$

Use Cramer's rule to find x_3 .

40. State and prove a version of Cramer's rule for solving the system $\vec{x}A = \vec{b}$, where $A \in M_n(R)$ is invertible, $\vec{b} \in R^n$ is a given row vector, and $\vec{x} \in R^n$ is an unknown row vector.
41. Let $A \in M_n(R)$ have $\det(A) \neq 0_R$. Use Cramer's rule and Laplace expansions to derive the explicit formula for the entries of A^{-1} .
42. Prove: if $A \in M_n(R)$ is invertible, then $\det(A^{-1}) = \det(A)^{-1}$ in R .
43. Matrices $A, B \in M_n(R)$ are called *similar* iff there exists an invertible $S \in M_n(R)$ with $B = S^{-1}AS$. Prove that similar matrices have the same trace, determinant, and characteristic polynomial.
44. For a field F , recall $GL_n(F)$ is the group of invertible matrices $A \in M_n(F)$. Prove that $SL_n(F) = \{A \in M_n(F) : \det(A) = 1_F\}$ is a normal subgroup of $GL_n(F)$. $SL_n(F)$ is called the *special linear group of degree n over F*.
45. Let $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ and $B = \begin{bmatrix} r & s \\ t & u \end{bmatrix}$. Expand the definition of $\det(AB)$ to get a sum of 8 terms, and show how terms cancel to produce a sum of 4 terms equal to $\det(A)\det(B)$. (This illustrates the proof in §5.13.)
46. Let $A, B \in M_n(R)$. Give direct proofs of the product formula $\det(AB) = \det(A)\det(B)$ for each of the following special cases: (a) A is diagonal; (b) A is triangular (cf. Exercise 21); (c) A is a *permutation matrix* (i.e., every row and column of A has exactly one 1, with all other entries zero); (d) A can be factored as $A = PLU$ for some $P, L, U \in M_n(R)$ with P a permutation matrix, L lower-triangular, and U upper-triangular.
47. Compute both sides of the Cauchy–Binet formula for the matrices

$$A = \begin{bmatrix} 2 & 1 & 0 & 3 \\ 1 & -1 & 4 & 3 \\ 5 & -3 & -2 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 3 & 1 & 2 \\ 0 & -2 & 3 \\ 1 & 1 & 2 \\ 3 & -3 & 1 \end{bmatrix}.$$

48. Use the Cauchy–Binet formula to prove the *Cauchy–Schwarz inequality*: for all real $x_1, \dots, x_n, y_1, \dots, y_n$,

$$\left| \sum_{i=1}^n x_i y_i \right| \leq \left(\sum_{i=1}^n x_i^2 \right)^{1/2} \cdot \left(\sum_{i=1}^n y_i^2 \right)^{1/2}.$$

49. Find the characteristic polynomial of each matrix in Exercise 28, and verify by a direct computation that the Cayley–Hamilton theorem holds for these matrices.
50. Illustrate the proof of the Cayley–Hamilton theorem given in §5.15 for a general matrix $A \in M_3(R)$ by computing the entries of the matrices C , zC^T , and $zC^T \text{adj}(C^T)$ appearing in that proof.
51. This exercise fills in some details in the proof given in §5.15. (a) Show that the set S of matrices defined in that proof is a subring of $M_n(R)$, and S is commutative. (b) Show that $h : R \rightarrow S$, defined by $h(r) = rI_n$ for $r \in R$, is a ring homomorphism. (c) Prove that $(zU)V = z(UV)$ for all $U, V \in M_n(S)$ and all $z \in (R^n)^n$.
52. Compute the permanent of each matrix in Exercise 28.

53. Let $A \in M_n(R)$. Give detailed proofs of the following facts about permanents:
 (a) $\text{per}(A) = \text{per}(A^T)$; (b) $\text{per}(A^*) = \text{per}(A)$ for $R = \mathbb{C}$; (c) $\text{per} : (R^n)^n \rightarrow R$ is an R -multilinear function of the n rows of A (or the n columns of A); (d) if A has a row or column of zeroes, then $\text{per}(A) = 0$; (e) if A is triangular, then $\text{per}(A) = \prod_{i=1}^n A(i, i)$; (f) for $c \in R$, $\text{per}(cA) = c^n \text{per}(A)$; (g) for $i \in [n]$, $\text{per}(A) = \sum_{k=1}^n A(i, k) \text{per}(A[i|k])$.
54. Give an example with $A, B \in M_2(\mathbb{R})$ to show that $\text{per}(AB) = \text{per}(A)\text{per}(B)$ is false in general.
55. *Pfaffians*. For all even $n > 0$, let SPf_n be the set of $f \in S_n$ with $f(1) < f(3) < f(5) < \dots < f(n-1)$ and $f(i) < f(i+1)$ for all odd $i \in [n]$. Given $A \in M_n(R)$ with $A^T = -A$, define the *Pfaffian* of A by the formula

$$\text{Pf}(A) = \sum_{f \in \text{SPf}_n} \text{sgn}(f) A(f(1), f(2)) A(f(3), f(4)) \cdots A(f(n-1), f(n)).$$

(a) Compute $\text{Pf}(A)$ for $n \in \{2, 4, 6\}$. (b) Prove that, for even $n > 0$, $\text{Pf}(A)$ is a sum of $1 \cdot 3 \cdot 5 \cdot \dots \cdot (n-1)$ terms. (c) For $i < j$ in $[n]$, let $A[[i, j]]$ be the matrix obtained from A by deleting row i , row j , column i , and column j . For even $n \geq 4$, prove

$$\text{Pf}(A) = \sum_{j=2}^n (-1)^j A(1, j) \text{Pf}(A[[1, j]]).$$

(d) Prove that $\det(A) = \text{Pf}(A)^2$ for $n \in \{2, 4, 6\}$. Can you prove it for all even $n > 0$? (For one solution, consult §12.12 of [36].)

56. Let R be the *non-commutative* ring $M_2(\mathbb{R})$. Use formula (5.1) to define the determinant of any $A \in M_n(R)$. Give specific examples to show that all but two statements below are *false*. Prove the two true statements. (a) For all $r, s, t, u \in R$, $\det \begin{bmatrix} r & s \\ t & u \end{bmatrix} = ru - st$. (b) For all $A \in M_n(R)$, $\det(A) = \det(A^T)$. (c) Switching two columns of $A \in M_n(R)$ replaces $\det(A)$ by $-\det(A)$. (d) Switching two rows of $A \in M_n(R)$ replaces $\det(A)$ by $-\det(A)$. (e) For all $A, B \in M_n(R)$, $\det(AB) = \det(A)\det(B)$. (f) For all $A \in M_n(R)$,

$$\det(A) = \sum_{i=1}^n (-1)^{i+1} A(i, 1) \det(A[i|1]).$$

(g) For all $A \in M_n(R)$,

$$\det(A) = \sum_{j=1}^n (-1)^{1+j} A(1, j) \det(A[1|j]).$$

57. Let $A \in M_{k,n}(R)$ with $k \leq n$. Fix lists (i_1, i_2, \dots, i_k) and (j_1, j_2, \dots, j_k) in $[n]^k$ and fix $s \in [k]$. Prove

$$\begin{aligned} & \det(A^{[i_1]}, A^{[i_2]}, \dots, A^{[i_k]}) \det(A^{[j_1]}, A^{[j_2]}, \dots, A^{[j_k]}) \\ &= \sum_{t=1}^k \det(A^{[i_1]}, \dots, A^{[j_t]}, \dots, A^{[i_k]}) \det(A^{[j_1]}, \dots, A^{[i_s]}, \dots, A^{[j_k]}), \end{aligned}$$

where the matrices on the right side arise from the matrices on the left side by interchanging the positions of columns $A^{[i_s]}$ and $A^{[j_t]}$.

58. True or false? Explain each answer. (a) For all $A, B \in M_n(R)$, $\det(A + B) = \det(A) + \det(B)$. (b) For all $A \in M_n(R)$ and $c \in R$, $\det(cA) = c\det(A)$. (c) For all $A \in M_n(\mathbb{C})$, if $A^k = I_n$ for some $k > 0$, then $|\det(A)| = 1$. (d) For all upper-triangular $A \in M_n(R)$ and all strictly upper-triangular $B \in M_n(R)$, $\det(A + B) = \det(A) + \det(B)$. (e) For all $n \times n$ matrices A with integer entries, if $\det(A) = \pm 1$, then every entry of A^{-1} is also an integer. (f) For all $A \in M_{m,n}(\mathbb{R})$ and $B \in M_{n,m}(\mathbb{R})$, if $m > n$ then $\det(AB) = 0$. (g) For all $A \in M_n(R)$, A^{-1} exists in $M_n(R)$ iff $\det(A) \neq 0_R$. (h) For all $A \in M_n(\mathbb{R})$, if $|A(i,j)| \leq K$ for all $i, j \in [n]$, then $|\det(A)| \leq n!K^n$. (i) For all $n \in \mathbb{N}^+$ and all $A, B \in M_n(\mathbb{Z}_2)$, $\text{per}(AB) = \text{per}(A)\text{per}(B)$.

This page intentionally left blank

6

Concrete vs. Abstract Linear Algebra

In elementary linear algebra, we learn about column vectors, matrices, and the algebraic operations on these objects: vector addition, multiplication of a vector by a scalar, matrix addition, multiplication of a matrix by a scalar, matrix multiplication, matrix inversion, matrix transpose, etc. Later in linear algebra, we study abstract vector spaces and linear maps between such spaces. This abstract setting provides a powerful tool for theoretical work. But for computational applications, it is often more convenient to work with column vectors and matrices. The goal of this chapter is to give a thorough explanation of the relation between the concrete world of column vectors and matrices on the one hand, and the abstract world of vector spaces and linear maps on the other hand. We will build a dictionary linking these two worlds, one entry at a time, which explains the exact connection between each abstract concept and its concrete counterpart. Our complete matrix-theory/linear-algebra dictionary appears in the summary for this chapter.

We will mainly be interested in the following three abstract entities: vectors, linear transformations, and bilinear maps. Given a vector v from an abstract n -dimensional vector space V over a field F , we will see how to represent v by a concrete column vector (or n -tuple) with entries from F . Given a linear transformation $T : V \rightarrow W$ mapping an abstract n -dimensional F -vector space V to an abstract m -dimensional F -vector space W , we will see how to represent T by a concrete $m \times n$ matrix with entries from F . Finally, given a bilinear form B on an abstract n -dimensional vector space V , we will see how to represent B by a concrete $n \times n$ matrix with entries in F . For each of these constructions, we will discuss how the explicit concrete operations on column vectors and matrices correspond to abstract algebraic operations on vectors, linear maps, and bilinear maps.

All of our constructions for converting abstract entities to concrete entities depend heavily on choosing ordered bases for the abstract vector spaces in question. We will see that different choices of bases lead to different concrete representations of a given vector, linear map, or bilinear form. One may then ask if some bases are better than others for computational or theoretical purposes. For instance, since it is easier to compute with diagonal matrices compared to arbitrary matrices, one could ask if there are bases such that a given linear operator is represented by a diagonal matrix. This leads us to a discussion of similarity, congruence, transition matrices, change of coordinates, diagonalization, and triangularization. The end of the chapter studies real and complex inner product spaces, orthogonal and unitary maps and matrices, and orthonormal bases.

6.1 Concrete Column Vectors vs. Abstract Vectors

In this chapter, we assume familiarity with the material on matrices covered in Chapter 4. For each field F and each positive integer n , we define a concrete vector space F^n whose elements are n -tuples (c_1, c_2, \dots, c_n) with each $c_i \in F$. For example, if $F = \mathbb{C}$ and $n = 3$,

some specific elements of \mathbb{C}^3 are $(i, \pi, e - 3i)$, $(7/2, 0, -i\sqrt{3})$, and $(0, 0, 1)$. We often identify

the n -tuple (c_1, c_2, \dots, c_n) with the column vector $\begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix}$.

On the other hand, let V be an abstract finite-dimensional vector space over F . In the absence of more specific information about V , it is not clear at the outset how to represent or describe particular elements in V . The best we can do is to introduce generic letters like v, w, \dots , to denote vectors in V .

The concept of *linear combinations* allows us to relate the concrete set F^n to the abstract set V . Suppose $X = (x_1, \dots, x_n)$ is any ordered list of vectors in V . Define a *linear combination function* $L_X : F^n \rightarrow V$ by setting

$$L_X(c_1, \dots, c_n) = c_1x_1 + c_2x_2 + \cdots + c_nx_n \in V \quad \text{for all } (c_1, \dots, c_n) \in F^n.$$

The function L_X maps an n -tuple of scalars (c_1, \dots, c_n) to the linear combination of the x_i 's with the c_i 's as coefficients.

We may now seek conditions under which the function L_X will be surjective, injective, or bijective. First, the following conditions are logically equivalent: L_X is surjective; for every $v \in V$, there exists $w \in F^n$ with $v = L_X(w)$; for all $v \in V$, there exist $c_1, \dots, c_n \in F$ with $v = \sum_{i=1}^n c_i x_i$; every $v \in V$ can be written in *at least one way* as a linear combination of the vectors in X ; the list X *spans* V .

Second, the following conditions are logically equivalent: L_X is injective; for all $w, y \in F^n$, $L_X(w) = L_X(y)$ implies $w = y$; for all $c_1, \dots, c_n, d_1, \dots, d_n \in F$, $\sum_{i=1}^n c_i x_i = \sum_{i=1}^n d_i x_i$ implies $(c_1, \dots, c_n) = (d_1, \dots, d_n)$; each $v \in V$ can be written in *at most one way* as a linear combination of the vectors in X ; the list X is *linearly independent*. To see why linear independence of X is equivalent to the preceding conditions, assume that X is linearly independent. Given $c_i, d_i \in F$ with $\sum_{i=1}^n c_i x_i = \sum_{i=1}^n d_i x_i$, note that $\sum_{i=1}^n (c_i - d_i)x_i = 0$. By the linear independence of X , we conclude $c_i - d_i = 0$ and $c_i = d_i$ for all i . Conversely, suppose L_X is injective and $\sum_{i=1}^n c_i x_i = 0$ for some $c_i \in F$. Then $L_X(c_1, \dots, c_n) = 0_V = L_X(0, \dots, 0)$ forces $(c_1, \dots, c_n) = (0, \dots, 0)$, so all $c_i = 0$. This means that the list X is linearly independent.

Combining the remarks in the last two paragraphs, we see that the following conditions are equivalent: $L_X : F^n \rightarrow V$ is a bijection; L_X is surjective and injective; X spans V and X is linearly independent; X is an *ordered basis* of V ; every $v \in V$ can be written in *exactly one way* in the form $v = \sum_{i=1}^n c_i x_i = L_X(c_1, \dots, c_n)$ for some $c_i \in F$.

Assume henceforth that $X = (x_1, \dots, x_n)$ is an ordered basis of V , so that L_X is bijective and hence invertible. For all $v \in V$, define $[v]_X = L_X^{-1}(v) \in F^n$. So, for all $c_i \in F$,

$$[v]_X = (c_1, \dots, c_n) \text{ iff } v = c_1x_1 + c_2x_2 + \cdots + c_nx_n,$$

and for each $v \in V$ there exist unique scalars $c_1, \dots, c_n \in F$ for which these equalities will hold. We call the scalars c_1, \dots, c_n the *coordinates of v relative to the ordered basis X* , and we call $[v]_X$ the *coordinate vector of v relative to X* . For example, fix $j \in [n] = \{1, 2, \dots, n\}$. Let $e_j \in F^n$ be the n -tuple $(0, \dots, 1, \dots, 0)$ that has a 1 in position j and zeroes elsewhere. Since $L_X(e_j) = 0x_1 + \cdots + 1x_j + \cdots + 0x_n = x_j$, we have $[x_j]_X = e_j$.

To summarize: for every ordered basis X of V , there is an associated bijection L_X from the concrete vector space F^n to the abstract vector space V , which maps an n -tuple to the corresponding linear combination of elements of X . The inverse of L_X , denoted $v \mapsto [v]_X$, lets us describe abstract vectors by concrete n -tuples (coordinates). For $c_i \in F$ and $v \in V$, $[v]_X = (c_1, \dots, c_n)$ iff $v = c_1x_1 + \cdots + c_nx_n$. As suggested by the notation, the coordinate

vector $[v]_X \in F^n$ depends on both the vector v and the choice of ordered basis X . Changing to a different ordered basis Y will replace $[v]_X$ with a new n -tuple $[v]_Y$. The relationship between $[v]_X$ and $[v]_Y$ will be explained in §6.10.

6.2 Examples of Computing Coordinates

Throughout this chapter, we will use the real vector spaces and ordered bases shown in Table 6.1 as running examples of the computations discussed in the text. We let the reader check that each list X_i displayed in the table really is an ordered basis for the given vector space.

TABLE 6.1

Vector Spaces and Ordered Bases Used as Examples in This Chapter.

1. $\mathbb{R}^3 = \{(c_1, c_2, c_3) : c_i \in \mathbb{R}\}$ (3-tuples or column vectors). Basis $X_1 = (e_1, e_2, e_3) = ((1, 0, 0), (0, 1, 0), (0, 0, 1))$. Basis $X_2 = ((1, 1, 1), (1, 2, 4), (1, 3, 9))$. Basis $X_3 = ((0, 2, 1), (-1, 1, 3), (1, 0, 4))$.
2. $P_{\leq 3} = \{f \in \mathbb{R}[t] : f = 0 \text{ or } \deg(f) \leq 3\}$ (polynomials in t of degree at most 3). Basis $X_1 = (1, t, t^2, t^3)$. Basis $X_2 = ((t - 2)^3, (t - 2)^2, (t - 2), 1)$. Basis $X_3 = (t + 3, 2t^2 - 4, t^3 - t^2, t^3 + t^2)$.
3. $M_2(\mathbb{R}) = \left\{ \begin{bmatrix} a & b \\ c & d \end{bmatrix} : a, b, c, d \in \mathbb{R} \right\}$ (2×2 matrices). Basis $X_1 = (e_{11}, e_{12}, e_{21}, e_{22}) = \left(\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \right)$. Basis $X_2 = \left(\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}, \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \right)$. Basis $X_3 = \left(\begin{bmatrix} 1 & 2 \\ 0 & 3 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 2 & 3 \end{bmatrix}, \begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix}, \begin{bmatrix} 1 & 2 \\ 3 & 0 \end{bmatrix} \right)$.
4. $\mathbb{C} = \{a + ib : a, b \in \mathbb{R}\}$ (complex numbers). Basis $X_1 = (1, i)$. Basis $X_2 = (e^{\pi i/6}, e^{2\pi i/3})$. Basis $X_3 = (3 + 4i, 2 - i)$.

To **compute the coordinate vector** $[v]_X$ given $v \in V$ and an ordered basis $X = (x_1, \dots, x_n)$: introduce unknown scalars $c_1, \dots, c_n \in F$; solve the equation $v = c_1x_1 + \dots + c_nx_n$ for these scalars; output the answer $[v]_X = (c_1, \dots, c_n)$. Typically, the vector equation $v = \sum_{i=1}^n c_i x_i$ can be rewritten as a system of n linear scalar equations in n unknowns c_i , which can be solved by Gaussian elimination or similar methods.

Example 1. Let $v = (3, 3, 1) \in \mathbb{R}^3$. Since $v = 3(1, 0, 0) + 3(0, 1, 0) + 1(0, 0, 1)$, we see that $[v]_{X_1} = (3, 3, 1) = v$. (More generally, for any $w \in F^n$, $[w]_X = w$ for the standard

ordered basis $X = (e_1, \dots, e_n)$.) On the other hand, to find $[v]_{X_2}$, we must solve the vector equation

$$(3, 3, 1) = c_1(1, 1, 1) + c_2(1, 2, 4) + c_3(1, 3, 9).$$

Equating components gives a system of three linear equations

$$3 = 1c_1 + 1c_2 + 1c_3, \quad 3 = 1c_1 + 2c_2 + 3c_3, \quad 1 = 1c_1 + 4c_2 + 9c_3.$$

Solving this system gives $c_1 = 2$, $c_2 = 2$, and $c_3 = -1$. Therefore $[v]_{X_2} = (2, 2, -1)$. To find coordinates relative to X_3 , we must now solve

$$(3, 3, 1) = c_1(0, 2, 1) + c_2(-1, 1, 3) + c_3(1, 0, 4),$$

or equivalently

$$3 = 0c_1 - 1c_2 + 1c_3, \quad 3 = 2c_1 + 1c_2 + 0c_3, \quad 1 = 1c_1 + 3c_2 + 4c_3.$$

The solution is $[v]_{X_3} = (32/13, -25/13, 14/13)$.

Example 2. Let $v = (t-3)^3 \in P_{\leq 3}$. To find $[v]_{X_1}$, solve $(t-3)^3 = c_11 + c_2t + c_3t^2 + c_4t^3$ by expanding the left side into monomials and comparing coefficients. Since $(t-3)^3 = t^3 - 9t^2 + 27t - 27$, we get $c_1 = -27$, $c_2 = 27$, $c_3 = -9$, and $c_4 = 1$, so $[v]_{X_1} = (-27, 27, -9, 1)$. On the other hand, to find $[v]_{X_2}$, we must solve

$$(t-3)^3 = d_1(t-2)^3 + d_2(t-2)^2 + d_3(t-2) + d_41$$

for unknowns $d_1, d_2, d_3, d_4 \in \mathbb{R}$. Expanding both sides gives

$$t^3 - 9t^2 + 27t - 27 = d_1t^3 + (d_2 - 6d_1)t^2 + (d_3 - 4d_2 + 12d_1)t + (d_4 - 2d_3 + 4d_2 - 8d_1).$$

Equating coefficients leads to the system of linear equations

$$d_1 = 1, \quad d_2 - 6d_1 = -9, \quad d_3 - 4d_2 + 12d_1 = 27, \quad d_4 - 2d_3 + 4d_2 - 8d_1 = -27,$$

which has solution $d_1 = 1$, $d_2 = -3$, $d_3 = 3$, $d_4 = -1$. So $[v]_{X_2} = (1, -3, 3, -1)$. Similarly, solving more linear equations leads to $[v]_{X_3} = (27, 27, 32, -31)$.

Example 3. Let $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \in M_2(\mathbb{R})$. We see by inspection that $A = 1e_{11} + 2e_{12} + 3e_{21} + 4e_{22}$, so that $[A]_{X_1} = (1, 2, 3, 4)$. On the other hand, writing A as a linear combination of vectors in X_2 using unknown scalars d_i , and equating corresponding matrix entries on each side, we obtain the linear equations

$$d_1 + d_3 = 1, \quad d_2 - d_4 = 2, \quad d_2 + d_4 = 3, \quad d_1 - d_3 = 4,$$

which have solution $[A]_{X_2} = (d_1, d_2, d_3, d_4) = (5/2, 5/2, -3/2, 1/2)$. Similarly, by solving

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = c_1 \begin{bmatrix} 1 & 2 \\ 0 & 3 \end{bmatrix} + c_2 \begin{bmatrix} 1 & 0 \\ 2 & 3 \end{bmatrix} + c_3 \begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix} + c_4 \begin{bmatrix} 1 & 2 \\ 3 & 0 \end{bmatrix},$$

we obtain $[A]_{X_3} = (1/3, 1/3, 2/3, 1/3)$.

Example 4. Let $v = e^{\pi i/3} \in \mathbb{C}$. By definition of complex exponentials, we have $v = \cos(\pi/3)1 + \sin(\pi/3)i = (1/2)1 + (\sqrt{3}/2)i$, and hence $[v]_{X_1} = (1/2, \sqrt{3}/2)$. To find $[v]_{X_2}$, we write $e^{\pi i/3} = d_1e^{\pi i/6} + d_2e^{2\pi i/3}$ for real unknowns d_1 and d_2 . Equating real and imaginary parts of both sides leads to the two linear equations

$$1/2 = (\sqrt{3}/2)d_1 + (-1/2)d_2, \quad \sqrt{3}/2 = (1/2)d_1 + (\sqrt{3}/2)d_2.$$

The solution is $d_1 = \sqrt{3}/2$, $d_2 = 1/2$, so $[v]_{X_2} = (\sqrt{3}/2, 1/2)$. We note in passing that if we reverse the order of the basis X_1 , we get $[v]_{(i,1)} = (\sqrt{3}/2, 1/2) = [v]_{X_2}$, showing that a vector can have the same coordinates relative to two different ordered bases. Finally, solving two more linear equations shows that $[v]_{X_3} = ((1+2\sqrt{3})/22, (4-3\sqrt{3})/22)$.

6.3 Concrete vs. Abstract Vector Space Operations

Recall that vector addition and scalar multiplication in F^n are defined by the explicit formulas

$$(c_1, \dots, c_n) + (d_1, \dots, d_n) = (c_1 + d_1, \dots, c_n + d_n), \quad a(c_1, \dots, c_n) = (ac_1, \dots, ac_n)$$

for all $c_i, d_i, a \in F$. On the other hand, all we know initially about the addition and scalar multiplication operations in an abstract vector space V are the axioms listed in Table 1.4.

We now study how the linear combination maps $L_X : F^n \rightarrow V$ relate the concrete vector space operations in F^n to the abstract vector space operations in V . First, let $X = (x_1, \dots, x_n)$ be any ordered list in V . For $\mathbf{c} = (c_1, \dots, c_n) \in F^n$, $\mathbf{d} = (d_1, \dots, d_n) \in F^n$, and $a \in F$, we use the vector space axioms in V to compute:

$$\begin{aligned} L_X(\mathbf{c} + \mathbf{d}) &= L_X(c_1 + d_1, \dots, c_n + d_n) = \sum_{i=1}^n (c_i + d_i)x_i = \sum_{i=1}^n c_i x_i + \sum_{i=1}^n d_i x_i = L_X(\mathbf{c}) + L_X(\mathbf{d}); \\ L_X(a\mathbf{c}) &= L_X(ac_1, \dots, ac_n) = \sum_{i=1}^n (ac_i)x_i = a \sum_{i=1}^n c_i x_i = aL_X(\mathbf{c}). \end{aligned}$$

This computation shows that $L_X : F^n \rightarrow V$ is always an *F -linear transformation* from F^n into V . Combining this fact with the results in §6.1, we see that L_X is a vector space isomorphism from F^n to V iff X is an ordered basis for V . In this case, the map $L_X^{-1} : V \rightarrow F^n$, given by $v \mapsto [v]_X$ for $v \in V$, is also a vector space isomorphism (see Exercise 38 of Chapter 1). Linearity of this inverse map means that

$$[v + w]_X = [v]_X + [w]_X; \quad [cv]_X = c[v]_X \quad \text{for all } v, w \in V \text{ and } c \in F.$$

Example. Let us find $[e_{22}]_{X_3}$ using computations from §6.2. Let $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ and $B = \begin{bmatrix} 1 & 2 \\ 3 & 0 \end{bmatrix}$, and note $e_{22} = (1/4)A - (1/4)B$. We already calculated $[A]_{X_3} = (1/3, 1/3, 2/3, 1/3)$, and evidently $[B]_{X_3} = (0, 0, 0, 1)$. Therefore,

$$[e_{22}]_{X_3} = [(1/4)(A - B)]_{X_3} = (1/4)[A]_{X_3} - (1/4)[B]_{X_3} = (1/12, 1/12, 1/6, -1/6).$$

To summarize: each ordered basis X of V gives us a vector space isomorphism from the abstract space V to the concrete space F^n (where $n = \dim(V)$ is the length of X), given by the coordinate map $v \mapsto [v]_X$ for $v \in V$. The inverse isomorphism L_X sends $(c_1, \dots, c_n) \in F^n$ to $\sum_{i=1}^n c_i x_i \in V$. We have $L_X(e_j) = x_j$ and $[x_j]_X = e_j$, so that the standard ordered basis $E = (e_1, \dots, e_n)$ of F^n corresponds to the ordered basis X of V under the isomorphism. Note that we have found not one, but many isomorphisms $V \cong F^n$, which are indexed by the ordered bases of V . Many problems in linear algebra involve finding a “good” ordered basis X for V , so that the associated isomorphism $V \cong F^n$ has particularly nice properties.

Using the fact that any finite-dimensional vector space V has a unique dimension (which is the length of any ordered basis of V), we can now show that *two finite-dimensional F -vector spaces are isomorphic iff they have the same dimension*. For, if V and W both have dimension n , we can compose the isomorphisms $V \cong F^n$ and $F^n \cong W$ to conclude that V and W are isomorphic. Conversely, if $V \cong W$, then any isomorphism from V to W is a bijection mapping each ordered basis of V to an ordered basis of W . Thus ordered bases

of V and W must have the same length, so $\dim(V) = \dim(W)$. (For infinite-dimensional vector spaces, one can show that $V \cong W$ iff $\dim(V) = \dim(W)$, where the dimensions are viewed as infinite cardinals.)

We close this section with a somewhat technical result. We have exhibited a family of isomorphisms $V \cong F^n$ (or equivalently, isomorphisms $F^n \cong V$) parameterized by ordered bases of V . One can ask if we have found *all* such isomorphisms, and if our parameterization via ordered bases is *unique*. To pose this question formally, note that we already have a well-defined function ϕ , given by the formula $\phi(X) = L_X$, that maps the set \mathcal{B} of all ordered bases of V into the set \mathcal{I} of all vector space isomorphisms from F^n onto V . We will show that $\phi : \mathcal{B} \rightarrow \mathcal{I}$ is a bijection, which means that each $L \in \mathcal{I}$ has the form $\phi(X) = L_X$ for a unique $X \in \mathcal{B}$. In other words, *every isomorphism $L : F^n \rightarrow V$ has the form L_X for a unique ordered basis X of V* .

First we prove ϕ is injective. Assume $X = (x_1, \dots, x_n)$ and $Y = (y_1, \dots, y_n)$ are ordered bases of V such that $\phi(X) = \phi(Y)$, i.e., $L_X = L_Y$. We must show that $X = Y$. This follows because $x_j = L_X(e_j) = L_Y(e_j) = y_j$ for $1 \leq j \leq n$. Now we show that ϕ is surjective. Let $L : F^n \rightarrow V$ be any vector space isomorphism; we must find an ordered basis X such that $L = L_X = \phi(X)$. We find X by applying L to the standard ordered basis $E = (e_1, \dots, e_n)$ of F^n . Formally, define $x_j = L(e_j)$ for $1 \leq j \leq n$, and define $X = (x_1, \dots, x_n)$. Note that X is an ordered basis of V since it is the image of the ordered basis E of F^n under an isomorphism L . To check that the functions $L, L_X : F^n \rightarrow V$ are equal, we check that they have the same effect on each element $\mathbf{c} = (c_1, \dots, c_n) \in F^n$. Since $\mathbf{c} = \sum_{j=1}^n c_j e_j$, linearity of L gives

$$L(\mathbf{c}) = \sum_{j=1}^n c_j L(e_j) = \sum_{j=1}^n c_j x_j = L_X(\mathbf{c}).$$

This completes the proof that ϕ is a bijection.

6.4 Matrices vs. Linear Maps

The next step in our comparison of concrete matrix theory to abstract linear algebra is to establish a correspondence between matrices and linear transformations, which will be absolutely essential for all that follows. For any field F , let $M_{m,n}(F)$ be the set of all $m \times n$ matrices with entries in F . For any F -vector spaces V and W , let $L(V, W)$ be the set of all F -linear maps T from V to W . A function $T : V \rightarrow W$ lies in $L(V, W)$ iff $T(v + v') = T(v) + T(v')$ and $T(cv) = cT(v)$ for all $v, v' \in V$ and all $c \in F$.

Now suppose $\dim(V) = n$ and $\dim(W) = m$. We will construct not one but many bijections from $L(V, W)$ to $M_{m,n}(F)$, which will be parameterized by ordered bases of V and W . Specifically, let $X = (x_1, \dots, x_n)$ be any ordered basis for V and $Y = (y_1, \dots, y_m)$ be any ordered basis for W . We will define a bijection $M_{X,Y} : L(V, W) \rightarrow M_{m,n}(F)$ that sends a linear transformation T to a certain matrix $M_{X,Y}(T)$, which will be denoted ${}_Y[T]_X$. The matrix ${}_Y[T]_X$ is called *the matrix of T relative to the input basis X and the output basis Y* .

We will build the correspondence between linear maps and matrices by composing three bijections

$$L(V, W) \xrightarrow{R_X} W^n \xrightarrow{(L_Y^{-1})^{\times n}} (F^m)^n \longrightarrow M_{m,n}(F), \text{ where:} \quad (6.1)$$

$$\begin{aligned}
L(V, W) &= \text{set of } F\text{-linear maps } T : V \rightarrow W; \\
W^n &= \text{set of lists } (w_1, \dots, w_n) \text{ with all } w_j \in W; \\
(F^m)^n &= \text{set of lists } (z_1, \dots, z_n) \text{ with all } z_j \in F^m; \\
M_{m,n}(F) &= \text{set of } m \times n \text{ matrices with entries in } F.
\end{aligned}$$

The idea motivating the definition of R_X is the fact that a linear map T with domain V is completely determined by the values of T on the basis vectors in X . So, we define a *restriction map* $R_X : L(V, W) \rightarrow W^n$ by setting $R_X(T) = (T(x_1), T(x_2), \dots, T(x_n)) \in W^n$ for all $T \in L(V, W)$.

Let us check that R_X is injective and surjective. First, suppose $S, T \in L(V, W)$ are two linear maps such that $R_X(S) = R_X(T)$. This means $S(x_j) = T(x_j)$ for $1 \leq j \leq n$. We prove $S = T$ by showing $S(v) = T(v)$ for all $v \in V$. Given $v \in V$, write $v = \sum_{j=1}^n c_j x_j$ for some $c_j \in F$. Then, use linearity of S and T to compute

$$S(v) = S\left(\sum_{j=1}^n c_j x_j\right) = \sum_{j=1}^n c_j S(x_j) = \sum_{j=1}^n c_j T(x_j) = T\left(\sum_{j=1}^n c_j x_j\right) = T(v).$$

Next, given any $\mathbf{w} = (w_1, w_2, \dots, w_n) \in W^n$, we must build $T \in L(V, W)$ such that $R_X(T) = \mathbf{w}$. We define

$$T\left(\sum_{j=1}^n c_j x_j\right) = \sum_{j=1}^n c_j w_j \quad \text{for all } c_j \in F. \quad (6.2)$$

It is routine to check that T is a well-defined linear map from V to W such that $T(x_j) = w_j$ for $1 \leq j \leq n$. Hence, $T \in L(V, W)$ and $R_X(T) = (T(x_1), \dots, T(x_n)) = \mathbf{w}$. We now know R_X is a bijection. Moreover, for $\mathbf{w} \in W^n$, the proof of surjectivity shows that $R_X^{-1}(\mathbf{w})$ is the map $T \in L(V, W)$ defined by (6.2).

Next we define a bijection from W^n to $(F^m)^n$. Starting with $\mathbf{w} = (w_1, \dots, w_n) \in W^n$, we apply the bijection $L_Y^{-1} : W \rightarrow F^m$ to each element of the list \mathbf{w} to obtain a list $([w_1]_Y, \dots, [w_n]_Y)$ of n column vectors in F^m . The inverse bijection takes a list $(z_1, \dots, z_n) \in (F^m)^n$ and maps it to the list $(L_Y(z_1), \dots, L_Y(z_n)) \in W^n$.

Finally, as discussed in §4.1, there is a bijection from $(F^m)^n$ to $M_{m,n}(F)$ that maps a list (z_1, \dots, z_n) of n column vectors in F^m to the $m \times n$ matrix A whose columns are z_1, \dots, z_n in this order. The inverse map sends any $A \in M_{m,n}(F)$ to the list $(A^{[1]}, \dots, A^{[n]}) \in (F^m)^n$, where $A^{[j]} \in F^m$ denotes the j 'th column of A .

Composing the three bijections discussed above, we get a bijection $M_{X,Y}$ from $L(V, W)$ to $M_{m,n}(F)$. Applying the formulas for each map, we see that $T \in L(V, W)$ is first sent to the list $(T(x_1), \dots, T(x_n)) \in W^n$, which is then sent to the list of coordinate vectors $([T(x_1)]_Y, \dots, [T(x_n)]_Y) \in (F^m)^n$, which is finally sent to the matrix having these coordinate vectors as columns. So, writing ${}_Y[T]_X = M_{X,Y}(T)$, the j 'th column of this matrix is $({}_Y[T]_X)^{[j]} = [T(x_j)]_Y$ for $1 \leq j \leq n$. The inverse bijection starts with a matrix $A \in M_{m,n}(F)$, sends this matrix to its list of columns $(A^{[1]}, \dots, A^{[n]}) \in (F^m)^n$, then sends this list to the list of vectors $(w_1, \dots, w_n) \in W^n$, where $w_j = L_Y(A^{[j]}) = \sum_{i=1}^m A(i, j)y_i$ for $1 \leq j \leq n$, and finally sends this list to the linear map T defined in (6.2). Explicitly,

$$T\left(\sum_{j=1}^n c_j x_j\right) = \sum_{j=1}^n c_j \sum_{i=1}^m A(i, j)y_i = \sum_{i=1}^m \left(\sum_{j=1}^n A(i, j)c_j \right) y_i \quad \text{for all } c_j \in F. \quad (6.3)$$

To summarize: whenever $X = (x_1, \dots, x_n)$ is an ordered basis for V and $Y = (y_1, \dots, y_m)$ is an ordered basis for W , we have a bijection $T \mapsto {}_Y[T]_X$ from $L(V, W)$ to

$M_{m,n}(F)$. By definition, ${}_Y[T]_X$ is the matrix whose j 'th column consists of the coordinates of $T(x_j)$ relative to Y . A linear map T corresponds to a matrix $A = {}_Y[T]_X$ under this bijection iff $T(x_j) = \sum_{i=1}^m A(i,j)y_i$ for all $j \in [n]$ iff T is given by formula (6.3). As indicated by the notation, the bijections constructed here depend on the choice of ordered bases X and Y .

6.5 Examples of Matrices Associated with Linear Maps

We now give some examples of computing the matrices associated with linear transformations. **To compute the matrix** ${}_Y[T]_X$ of the linear map $T : V \rightarrow W$ relative to ordered bases $X = (x_1, \dots, x_n)$ of V and $Y = (y_1, \dots, y_m)$ of W : apply the linear map T to each basis vector $x_j \in X$; compute the coordinates of the resulting vector $T(x_j)$ relative to the basis Y ; and write down these coordinates as the j 'th column of the matrix.

Example 1. (See Table 6.1 for the notation used in these examples.) Let $T : P_{\leq 3} \rightarrow P_{\leq 3}$ be the differentiation operator $T(f) = df/dt$ for $f \in P_{\leq 3}$, which is known to be a linear map. Let us find $A = {}_{X_1}[T]_{X_1}$. To get the first column of A , compute $T(1) = 0$, which has coordinates $(0, 0, 0, 0)$ relative to X_1 . To get the second column of A , compute $T(t) = 1$, which has coordinates $(1, 0, 0, 0)$ relative to X_1 . To get the third column of A , compute $T(t^2) = 2t$, which has coordinates $(0, 2, 0, 0)$ relative to X_1 . To get the fourth column of A , compute $T(t^3) = 3t^2$, which has coordinates $(0, 0, 3, 0)$ relative to X_1 . Arranging these columns in a matrix, we have

$$A = {}_{X_1}[T]_{X_1} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 3 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

Now suppose we change the input basis to X_2 , but keep X_1 as the output basis. Applying T to the vectors $(t-2)^3, (t-2)^2, (t-2)^1, 1$ produces the vectors $3(t-2)^2 = 3t^2 - 12t + 12$, $2(t-2) = 2t - 4$, 1 , and 0 , respectively. We have $[3t^2 - 12t + 12]_{X_1} = (12, -12, 3, 0)$, and so on, leading to the matrix

$${}_{X_2}[T]_{X_2} = \begin{bmatrix} 12 & -4 & 1 & 0 \\ -12 & 2 & 0 & 0 \\ 3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

On the other hand, taking X_2 as both input basis and output basis gives

$${}_{X_2}[T]_{X_2} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 3 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}.$$

Example 2. Let $T : M_2(\mathbb{R}) \rightarrow M_2(\mathbb{R})$ be the transpose map given by $T(A) = A^T$ for all $A \in M_2(\mathbb{R})$; this map is linear. We compute $T(e_{11}) = e_{11}$, $T(e_{12}) = e_{21}$, $T(e_{21}) = e_{12}$, and $T(e_{22}) = e_{22}$. Therefore,

$${}_{X_1}[T]_{X_1} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

Similar computations show that

$$\begin{aligned} x_2[T]_{X_2} &= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1 \end{bmatrix}; \quad x_3[T]_{X_3} = \begin{bmatrix} 0 & 1 & 1/2 & 1/2 \\ 1 & 0 & -1/2 & -1/2 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}; \\ x_1[T]_{X_3} &= \begin{bmatrix} 1 & 1 & 0 & 1 \\ 0 & 2 & 2 & 3 \\ 2 & 0 & 1 & 2 \\ 3 & 3 & 3 & 0 \end{bmatrix}; \quad x_2[T]_{X_1} = \begin{bmatrix} 1/2 & 0 & 0 & 1/2 \\ 0 & 1/2 & 1/2 & 0 \\ 1/2 & 0 & 0 & -1/2 \\ 0 & 1/2 & -1/2 & 0 \end{bmatrix}. \end{aligned}$$

Example 3. For any vector spaces V and W , the zero map $0 \in L(V, W)$ is the linear map sending every $v \in V$ to 0_W . For all ordered bases X and Y , ${}_Y[0]_X$ is the zero matrix $0 \in M_{m,n}(F)$. On the other hand, for the identity map id_V on any vector space V , we have ${}_X[\text{id}]_X = I_n$ (the identity matrix) for all ordered bases X of V . But if X and Y are distinct ordered bases for V , then ${}_Y[\text{id}]_X$ is *not* the identity matrix. For example, using the bases

of \mathbb{R}^3 in Table 6.1, we see that ${}_X[\text{id}]_{X_2} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 4 & 9 \end{bmatrix}$. We study matrices of the form

${}_Y[\text{id}]_X$ in more detail later (see §6.10).

Example 4. Let $z = a+ib$ be a fixed complex number. Define $T_z : \mathbb{C} \rightarrow \mathbb{C}$ by $T_z(w) = zw$ for all $w \in \mathbb{C}$. One readily checks that T_z is an \mathbb{R} -linear map (and also a \mathbb{C} -linear map). Let us compute ${}_X[T_z]_{X_1}$. For column 1, compute $T_z(1) = z1 = z = a+ib$, and note that $[z]_{X_1} = (a, b)$. For column 2, compute $T_z(i) = zi = -b+ia$, and note that $[zi]_{X_1} = (-b, a)$. So

$${}_X[T_z]_{X_1} = \begin{bmatrix} a & -b \\ b & a \end{bmatrix}.$$

Example 5. Define a linear map $T : P_{\leq 3} \rightarrow \mathbb{R}^3$ by $T(f) = (f(1), f(2), f(3))$ for $f \in P_{\leq 3}$. Let $X = (1, t, t^2, t^3)$, $Y = (e_1, e_2, e_3)$, and $Z = ((1, 1, 1), (1, 2, 3), (1, 4, 9))$. Applying T to each vector in X , we get $T(1) = (1, 1, 1)$, $T(t) = (1, 2, 3)$, $T(t^2) = (1, 4, 9)$, and $T(t^3) = (1, 8, 27)$. Computing coordinates of these vectors relative to Y and Z , we find:

$${}_Y[T]_X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 4 & 8 \\ 1 & 3 & 9 & 27 \end{bmatrix}; \quad {}_Z[T]_X = \begin{bmatrix} 1 & 0 & 0 & 6 \\ 0 & 1 & 0 & -11 \\ 0 & 0 & 1 & 6 \end{bmatrix}.$$

If we change the input basis to $X' = (1, (t-2), (t-2)^2, (t-2)^3)$, we compute:

$${}_Y[T]_{X'} = \begin{bmatrix} 1 & -1 & 1 & -1 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{bmatrix}; \quad {}_Z[T]_{X'} = \begin{bmatrix} 1 & -2 & 4 & -2 \\ 0 & 1 & -4 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}.$$

Example 6. Given $V = F^n$, $W = F^m$, and any matrix $A \in M_{m,n}(F)$, define a linear map $L^A : F^n \rightarrow F^m$ by setting $L^A(v) = Av \in F^m$ for all column vectors $v \in F^n$. The map L^A is called *left multiplication by the matrix A*. Let $E = (e_1, \dots, e_n)$ be the standard ordered basis of F^n , and let $E' = (e'_1, \dots, e'_m)$ be the standard ordered basis of F^m . We will show that ${}_{E'}[L^A]_E = A$. To get the j 'th column of ${}_{E'}[L^A]_E$, we apply L^A to the j 'th vector in E , obtaining $L^A(e_j) = Ae_j = A^{[j]}$. Next, we take the coordinates of this column vector relative to E' . But for any column vector $w \in F^m$, $[w]_{E'} = w$, so the coordinate vector is $A^{[j]}$. Thus, the j 'th column of ${}_{E'}[L^A]_E$ equals the j 'th column of A for all j , proving that ${}_{E'}[L^A]_E = A$.

We can restate the result of this example by saying that *the matrix of left multiplication by A (relative to the standard ordered bases) is A itself*. It follows that the map $A \mapsto L^A$ from $M_{m,n}(F)$ to $L(F^n, F^m)$ is the inverse of the general bijection $T \mapsto {}_{E'}[T]_E$ constructed earlier, so this map is also a bijection. Bijectivity of the map $A \mapsto L^A$ means that *every linear transformation from F^n to F^m is left multiplication by a unique $m \times n$ matrix*.

6.6 Vector Operations on Matrices and Linear Maps

Recall that $M_{m,n}(F)$ is a vector space under the following operations: for all $A, B \in M_{m,n}(F)$ and $c \in F$, $A + B$ is the matrix with i, j -entry $A(i, j) + B(i, j)$, while cA is the matrix with i, j -entry $cA(i, j)$. We can also describe these operations in terms of the columns of the matrices via the formulas $(A + B)^{[j]} = A^{[j]} + B^{[j]}$ and $(cA)^{[j]} = c(A^{[j]})$ (see §4.7). The operations on the right sides of the last two formulas are addition and scalar multiplication of column vectors in the vector space F^m .

Now suppose S and T are two linear maps in $L(V, W)$ and $c \in F$. Define $S + T : V \rightarrow W$ by $(S + T)(v) = S(v) + T(v)$ for all $v \in V$; define $cT : V \rightarrow W$ by $(cT)(v) = c(T(v))$ for all $v \in V$. The operations on the right sides of these definitions are the abstract operations in the given vector space W . One verifies that $S + T$ and cT are indeed F -linear maps, so that we have the closure properties $S + T \in L(V, W)$ and $cT \in L(V, W)$. Also, the zero function is in $L(V, W)$. It follows that $L(V, W)$ is a subspace of the vector space of all functions from V to W (see Exercise 5 in Chapter 4), so that $L(V, W)$ is itself a vector space under the operations defined here.

Assume now that $X = (x_1, \dots, x_n)$ is an ordered basis for V , and that $Y = (y_1, \dots, y_m)$ is an ordered basis for W . We will prove that the bijection $M_{X,Y} : L(V, W) \rightarrow M_{m,n}(F)$ given by $M_{X,Y}(T) = {}_Y[T]_X$ is F -linear, hence is a vector space isomorphism. Let $S, T \in L(V, W)$ and $c \in F$. Computing each column, for all $j \in [n]$ we have

$$\begin{aligned} ({}_Y[S + T]_X)^{[j]} &= [(S + T)(x_j)]_Y = [S(x_j) + T(x_j)]_Y = [S(x_j)]_Y + [T(x_j)]_Y \\ &= ({}_Y[S]_X)^{[j]} + ({}_Y[T]_X)^{[j]}, \end{aligned}$$

and therefore $M_{X,Y}(S + T) = M_{X,Y}(S) + M_{X,Y}(T)$. We also have

$$({}_Y[cT]_X)^{[j]} = [(cT)(x_j)]_Y = [c(T(x_j))]_Y = c \cdot [T(x_j)]_Y = c \cdot ({}_Y[T]_X)^{[j]},$$

and therefore $M_{X,Y}(cT) = cM_{X,Y}(T)$. **To summarize:** given any ordered bases X of V and Y of W , the map $T \mapsto {}_Y[T]_X$ is a vector space isomorphism $L(V, W) \cong M_{m,n}(F)$.

Since isomorphic vector spaces have the same dimension, we deduce $\dim(L(V, W)) = \dim(M_{m,n}(F)) = mn$. We can obtain an even stronger conclusion by recalling that vector space isomorphisms send bases to bases. We have seen (§4.3) that the set $E = \{e_{ij} : i \in [m], j \in [n]\}$ is a basis for $M_{m,n}(F)$. Note that e_{ij} can be described as the $m \times n$ matrix whose j 'th column is $e_{ij}^{[j]} = e_i \in F^m$, and whose k 'th column (for $k \neq j$) is $e_{ij}^{[k]} = 0 \in F^m$. The image of E under $M_{X,Y}^{-1}$ is a basis for $L(V, W)$. Let us describe the linear map $T_{ij} = M_{X,Y}^{-1}(e_{ij})$. Since ${}_Y[T_{ij}]_X = e_{ij}$ by definition, we see that $[T_{ij}(x_k)]_Y = e_{ij}^{[k]} = 0$ for all $k \neq j$, so that $T_{ij}(x_k) = L_Y(0) = 0$ for all $k \neq j$. On the other hand, $[T_{ij}(x_j)]_Y = e_{ij}^{[j]} = e_i$ implies $T_{ij}(x_j) = L_Y(e_i) = y_i$. Thus T_{ij} sends the basis vector x_j to y_i and all other basis vectors in X to zero. By linearity, $T_{ij} : V \rightarrow W$ must be given by the formula

$$T_{ij} \left(\sum_{k=1}^n c_k x_k \right) = c_j y_i \quad \text{for all } c_k \in F. \tag{6.4}$$

6.7 Matrix Transpose vs. Dual Maps

Given a matrix $A \in M_{m,n}(F)$, recall that the *transpose* of A is the matrix $A^T \in M_{n,m}(F)$ defined by $A^T(i,j) = A(j,i)$ for $1 \leq i \leq n$ and $1 \leq j \leq m$. To understand the significance of this operation in the world of abstract linear algebra, we must first enter into a brief digression on dual spaces. (Dual spaces are discussed at greater length in Chapter 13.)

Given an n -dimensional F -vector space V , the *dual space* V^* is defined to be $L(V,F)$, the vector space of all linear maps from V into the field F (viewed as a 1-dimensional F -vector space). Applying our general results on $L(V,W)$ with $W = F$, we see that $V^* = L(V,F)$ is isomorphic to the vector space $M_{1,n}(F)$ of $1 \times n$ matrices or *row vectors*. More specifically, let $X = (x_1, \dots, x_n)$ be any ordered basis of V . The one-element list $Z = (1_F)$ is an ordered basis for the F -vector space F , as one immediately checks. For $1 \leq j \leq n$, define linear maps $f_j \in V^* = L(V,F)$ by setting $f_j(x_j) = 1_F$, setting $f_j(x_k) = 0_F$ for $k \neq j$, and extending by linearity. Explicitly, $f_j(\sum_{k=1}^n c_k x_k) = c_j$ for all $c_1, \dots, c_n \in F$. Comparing this to the description of the maps T_{ij} at the end of the last section, we see that $X^* = (f_1, \dots, f_n)$ is the ordered basis for $V^* = L(V,F)$ that corresponds to the standard ordered basis for $M_{1,n}(F)$ under the isomorphism $M_{X,Z}^{-1}$. X^* is called the *dual basis* of X . Note that $\dim(V^*) = \dim(V) = n$ (in the finite-dimensional case we are considering).

Now, let W be an m -dimensional F -vector space with ordered basis $Y = (y_1, \dots, y_m)$, and let $S \in L(V,W)$ be a fixed linear map. Given any map $g \in W^* = L(W,F)$, note that $g \circ S : V \rightarrow F$ is a linear map from V to F , i.e., an element of V^* . Thus, S induces a map $S^* : W^* \rightarrow V^*$ given by $S^*(g) = g \circ S$ for $g \in W^*$. The map S^* is itself linear, since for all $g, h \in W^*$ and all $c \in F$,

$$S^*(g + h) = (g + h) \circ S = (g \circ S) + (h \circ S) = S^*(g) + S^*(h);$$

$$S^*(cg) = (cg) \circ S = c(g \circ S) = cS^*(g).$$

So $S^* \in L(W^*, V^*)$. Moreover, we have the dual bases Y^* and X^* for W^* and V^* , respectively. We will now prove that

$$x^*[S^*]_{Y^*} = [{}_Y[S]_X]^T.$$

That is, transposing the matrix of S (relative to the input basis X and output basis Y) gives the matrix of S^* (relative to the input basis Y^* and output basis X^*).

To prove this statement, set $A = {}_Y[S]_X$ and $B = x^*[S^*]_{Y^*}$; we must show that $B(i,j) = A(j,i)$ for all $i \in [n]$ and $j \in [m]$. Write $X^* = (f_1, \dots, f_n)$ and $Y^* = (g_1, \dots, g_m)$. By definition, the j 'th column of B is

$$B^{[j]} = [S^*(g_j)]_{X^*} = [g_j \circ S]_{X^*}.$$

To proceed, recall that $[g_j \circ S]_{X^*}$ is the unique n -tuple $(c_1, \dots, c_n) \in F^n$ such that

$$g_j \circ S = c_1 f_1 + \cdots + c_n f_n.$$

To find the i 'th coordinate c_i of $B^{[j]}$, evaluate both sides of the previous identity at x_i . The right side is c_i , while the left side is $g_j(S(x_i))$. Therefore, $B(i,j) = B^{[j]}(i) = c_i = g_j(S(x_i))$. On the other hand, since $A^{[i]} = [S(x_i)]_Y$, we have $S(x_i) = \sum_{k=1}^m A^{[i]}(k)y_k = \sum_{k=1}^m A(k,i)y_k$. Applying g_j to both sides gives $g_j(S(x_i)) = A(j,i)$. We conclude that $B(i,j) = A(j,i)$.

6.8 Matrix/Vector Multiplication vs. Evaluation of Maps

If A is an $m \times n$ matrix and $z \in F^n$ is an n -tuple (which we may regard as a column vector), we can form the matrix-vector product Az , which is an m -tuple (or column vector) whose i 'th component is $\sum_{j=1}^n A(i, j)z(j)$. As shown in Chapter 4, we can also write this formula as $Az = \sum_{j=1}^n z(j)A^{[j]}$, which expresses Az as a linear combination of the columns $A^{[j]}$ of A .

On the other hand, given a linear map $T : V \rightarrow W$ and an abstract vector $v \in V$, we can apply the map T to v to obtain another vector $w = T(v) \in W$. Let $X = (x_1, \dots, x_n)$ and $Y = (y_1, \dots, y_m)$ be ordered bases for V and W . We will prove that

$$[T(v)]_Y = {}_Y[T]_X [v]_X. \quad (6.5)$$

To state this formula in words: the coordinate vector of the output of T at input v (relative to Y) is found by multiplying the matrix of T (relative to X and Y) by the coordinate vector of v (relative to X). The proof consists of three steps. First, writing $[v]_X = (c_1, \dots, c_n)$, we know $v = \sum_{j=1}^n c_j x_j$ by definition of coordinates. Second, note that $T(v) = \sum_{j=1}^n c_j T(x_j)$ by linearity of T . Third, applying the formula at the end of the previous paragraph to $A = {}_Y[T]_X$ and $z = [v]_X$, we get

$${}_Y[T]_X [v]_X = \sum_{j=1}^n [v]_X(j) ({}_Y[T]_X)^{[j]} = \sum_{j=1}^n c_j [T(x_j)]_Y = \left[\sum_{j=1}^n c_j T(x_j) \right]_Y = [T(v)]_Y.$$

6.9 Matrix Multiplication vs. Composition of Linear Maps

For matrices $B \in M_{p,m}(F)$ and $A \in M_{m,n}(F)$, recall that the *matrix product* $BA \in M_{p,n}(F)$ is defined by

$$(BA)(i, j) = \sum_{k=1}^m B(i, k)A(k, j) \quad \text{for } i \in [p] \text{ and } j \in [n].$$

For linear maps $T \in L(V, W)$ and $S \in L(W, U)$, the *composition* $S \circ T \in L(V, U)$ is defined by

$$(S \circ T)(v) = S(T(v)) \quad \text{for all } v \in V;$$

one checks immediately that $S \circ T$ is indeed linear.

We now compare these two operations. Let $X = (x_1, \dots, x_n)$, $Y = (y_1, \dots, y_m)$, and $Z = (z_1, \dots, z_p)$ be ordered bases for V , W , and U , respectively. Given S and T as above, let $A = {}_Y[T]_X \in M_{m,n}(F)$, let $B = {}_Z[S]_Y \in M_{p,m}(F)$, and let $C = {}_Z[S \circ T]_X \in M_{p,n}(F)$. We will show that $C = BA$, i.e.,

$${}_Z[S \circ T]_X = {}_Z[S]_Y {}_Y[T]_X. \quad (6.6)$$

We saw in §4.7 that the j 'th column of BA is $B(A^{[j]})$, so it suffices to show that $C^{[j]} = B(A^{[j]})$ for $1 \leq j \leq n$. This follows from the computation:

$$C^{[j]} = [(S \circ T)(x_j)]_Z = [S(T(x_j))]_Z = {}_Z[S]_Y [T(x_j)]_Y = B(A^{[j]}),$$

where we have used the definition of C , then the definition of $S \circ T$, then (6.5), then the

definition of $A^{[j]}$. **To summarize:** the matrix of the composition of two linear maps is the product of the matrices of the individual maps, provided that the output basis of the map acting first matches the input basis of the map acting second.

As one consequence of this formula, suppose $U = V$, $m = n = p$, $X = Z$, and $T \in L(V, W)$ is an invertible linear map with inverse $T^{-1} \in L(W, V)$. Formula (6.6) gives

$${}_Y[T]_X \cdot {}_X[T^{-1}]_Y = {}_Y[T \circ T^{-1}]_Y = {}_Y[\text{id}_W]_Y = I_n;$$

$${}_X[T^{-1}]_Y \cdot {}_Y[T]_X = {}_X[T^{-1} \circ T]_X = {}_X[\text{id}_V]_X = I_n.$$

These equations show that

$$({}_Y[T]_X)^{-1} = {}_X[T^{-1}]_Y. \quad (6.7)$$

Thus, *matrix inversion corresponds to functional inversion of the associated linear map.*

6.10 Transition Matrices and Changing Coordinates

Given an element v of some vector space V , one often needs to compute the coordinates of v relative to more than one basis. Specifically, suppose X and Y are ordered bases for V , we are given $[v]_X$ for some $v \in V$, and we want to compute $[v]_Y$. Applying formula (6.5) to the identity map $\text{id} : V \rightarrow V$, we see that

$$[v]_Y = {}_Y[\text{id}]_X [v]_X.$$

Here, ${}_Y[\text{id}]_X$ is the matrix whose j 'th column is $[\text{id}(x_j)]_Y = [x_j]_Y$. We call this matrix the *transition matrix from X to Y* , since left-multiplication by this matrix transforms coordinates relative to X into coordinates relative to Y . Since $\text{id}^{-1} = \text{id}$, (6.7) shows that $({}_Y[\text{id}]_X)^{-1} = {}_X[\text{id}]_Y$. Thus, transition matrices between ordered bases are invertible, and the transition matrix from Y to X is the inverse of the transition matrix from X to Y .

To compute the transition matrix ${}_Y[\text{id}]_X$: compute the coordinates $[x_j]_Y$ of each x_j relative to Y and put these coordinates in the j 'th column of the transition matrix. **To change coordinates from $[v]_X$ to $[v]_Y$:** compute the matrix-vector product $[v]_Y = {}_Y[\text{id}]_X [v]_X$.

Example 1. For $V = P_{\leq 3}$ (see Table 6.1), we have

$${}_{X_1}[\text{id}]_{X_2} = \begin{bmatrix} -8 & 4 & -2 & 1 \\ 12 & -4 & 1 & 0 \\ -6 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}.$$

For instance, the first column arises from the computation $(t-2)^3 = -8 + 12t - 6t^2 + t^3$. The reader may check directly that

$${}_{X_2}[\text{id}]_{X_1} = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 6 \\ 0 & 1 & 4 & 12 \\ 1 & 2 & 4 & 8 \end{bmatrix},$$

and that this matrix is the inverse of the preceding one. Suppose $v = 2t^3 - 4t + 1$, so that $[v]_{X_1} = (1, -4, 0, 2)$. Then

$$[v]_{X_2} = {}_{X_2}[\text{id}]_{X_1} [v]_{X_1} = (2, 12, 20, 9),$$

which says that $v = 2(t - 2)^3 + 12(t - 2)^2 + 20(t - 2) + 9(1)$.

Example 2. For the three ordered bases of \mathbb{R}^3 in Table 6.1, it is immediate that

$$x_1[\text{id}]_{X_2} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 1 & 4 & 9 \end{bmatrix}, \quad x_1[\text{id}]_{X_3} = \begin{bmatrix} 0 & -1 & 1 \\ 2 & 1 & 0 \\ 1 & 3 & 4 \end{bmatrix}.$$

So, for instance,

$$x_2[\text{id}]_{X_3} = x_2[\text{id}]_{X_1} x_1[\text{id}]_{X_3} = (x_1[\text{id}]_{X_2})^{-1} x_1[\text{id}]_{X_3} = \begin{bmatrix} -4.5 & -4 & 5 \\ 7 & 4 & -7 \\ -2.5 & -1 & 3 \end{bmatrix}.$$

To double-check the third column of the answer, observe that $(1, 0, 4) = 5(1, 1, 1) - 7(1, 2, 4) + 3(1, 3, 9)$.

It turns out that all invertible matrices arise as transition matrices $Y[\text{id}]_X$. More precisely, suppose $\dim(V) = n$ and Y is a fixed ordered basis for V . Let \mathcal{B} be the set of all ordered bases for V , and let $\text{GL}_n(F)$ be the set of all invertible $n \times n$ matrices with entries in F . Define $\phi : \mathcal{B} \rightarrow \text{GL}_n(F)$ by setting $\phi(Z) = Y[\text{id}]_Z$ for every ordered basis $Z \in \mathcal{B}$. We will show that ϕ is a bijection.

For the proof, we define a map $\phi' : \text{GL}_n(F) \rightarrow \mathcal{B}$ that will be shown to be ϕ^{-1} . For any invertible matrix $A \in \text{GL}_n(F)$, let

$$\phi'(A) = (L_Y(A^{[1]}), \dots, L_Y(A^{[n]})).$$

We first show that $\phi'(A)$ is an ordered basis for V , i.e., $\phi'(A) \in \mathcal{B}$. Since A is invertible, the column rank of A is n (see §4.13). Consequently, the list of columns $(A^{[1]}, \dots, A^{[n]})$ is an ordered basis for F^n . Applying the isomorphism $L_Y : F^n \rightarrow V$ to this ordered basis, we see that $\phi'(A)$ is indeed an ordered basis for V . Now, let $A \in \text{GL}_n(F)$ and $Z = (z_1, \dots, z_n) \in \mathcal{B}$ be arbitrary. By definition of $Y[\text{id}]_Z$, we have

$$\phi(Z) = A \Leftrightarrow [z_j]_Y = A^{[j]} \text{ for all } j \in [n] \Leftrightarrow z_j = L_Y(A^{[j]}) \text{ for all } j \in [n] \Leftrightarrow Z = \phi'(A).$$

We conclude that ϕ' is the two-sided inverse to ϕ .

The bijectivity of ϕ means that *for any fixed ordered basis Y of V , every invertible matrix in $M_n(F)$ has the form $Y[\text{id}]_Z$ for a unique ordered basis Z of V* . One sees similarly that *for fixed Y , every invertible matrix in $M_n(F)$ has the form $X[\text{id}]_Y$ for a unique ordered basis X of V* .

6.11 Changing Bases

Let $\dim(V) = n$ and $\dim(W) = m$. Suppose $T \in L(V, W)$ is represented by a matrix $A = Y[T]_X$ relative to an ordered basis X for V and an ordered basis Y for W . What happens to A if we change the basis for the domain V from X to another ordered basis X' ? Since $T = T \circ \text{id}_V$, formula (6.6) gives the relation

$$Y[T]_{X'} = Y[T]_X X[\text{id}_V]_{X'}.$$

Writing $A' = Y[T]_{X'}$ and $P = X[\text{id}_V]_{X'}$, this says that $A' = AP$. Moreover, as X' ranges over all possible ordered bases for V (while X remains fixed), we know that P ranges over all the invertible $n \times n$ matrices.

For $A, A' \in M_{m,n}(F)$, say that A is *column-equivalent* to A' iff $A' = AP$ for some invertible $P \in M_n(F)$. One may check that this defines an equivalence relation on $M_{m,n}(F)$. Since P is invertible iff P is a finite product of elementary matrices, we see that A' is column-equivalent to A iff A' can be obtained from A by a finite sequence of elementary column operations, which are encoded by the matrix P (see Chapter 4). For a matrix $A = {}_Y[T]_X$, the equivalence class of A relative to column-equivalence is the set of all matrices ${}_{Y'}[T]_{X'}$, as X' ranges over all ordered bases for V .

Next, we ask: what happens to $A = {}_Y[T]_X$ if we change the basis for the codomain W from Y to Y' ? Since $T = \text{id}_W \circ T$, (6.6) gives

$${}_{Y'}[T]_X = {}_{Y'}[\text{id}_W]_Y {}_Y[T]_X.$$

Writing $A' = {}_{Y'}[T]_X$ and $Q = {}_{Y'}[\text{id}_W]_Y$, this says that $A' = QA$. Moreover, as Y' ranges over all possible ordered bases for W (while Y remains fixed), we know that Q ranges over all the invertible $m \times m$ matrices.

For $A, A' \in M_{m,n}(F)$, say that A is *row-equivalent* to A' iff $A' = QA$ for some invertible $Q \in M_m(F)$. This defines another equivalence relation on $M_{m,n}(F)$. A' is row-equivalent to A iff A' can be obtained from A by a finite sequence of elementary row operations, which are encoded by the matrix Q . If $A = {}_Y[T]_X$, then the equivalence class of A relative to row-equivalence is the set of all matrices ${}_{Y'}[T]_X$, as Y' ranges over all ordered bases for W .

Finally, if we change the input basis from X to X' and also change the output basis from Y to Y' , we have

$${}_{Y'}[T]_{X'} = {}_{Y'}[\text{id}_W]_Y {}_Y[T]_X {}_X[\text{id}_V]_{X'}.$$

Accordingly, we say that $A, A' \in M_{m,n}(F)$ are *equivalent* (or, more precisely, *row/column-equivalent*) iff $A' = QAP$ for some invertible matrices $Q \in M_m(F)$ and $P \in M_n(F)$ iff A' can be obtained from A by doing elementary row and/or column operations. The equivalence class of ${}_Y[T]_X$ relative to this equivalence relation is the set of all matrices ${}_{Y'}[T]_{X'}$ as X' and Y' range independently over all ordered bases for V and W , respectively.

If F is a field and $A \in M_{m,n}(F)$ is any matrix, we can find invertible matrices P and Q such that QAP has i, i -entry equal to 1 for $1 \leq i \leq r = \text{rank}(A)$, and all other entries of QAP are zero (see §4.12). Taking $A = {}_Y[T]_X$, where $T \in L(V, W)$ is an arbitrary linear map, we obtain the following *projection theorem*: for any $T \in L(V, W)$, there exist ordered bases $X' = (x'_1, \dots, x'_n)$ for V and $Y' = (y'_1, \dots, y'_m)$ for W such that $T(x'_i) = y'_i$ for $1 \leq i \leq \text{rank}(T)$, and $T(x'_i) = 0$ for all $i > \text{rank}(T)$. (One can also give an abstract proof of this theorem that makes no appeal to matrices; see Exercise 45.) In Chapter 18, we will obtain generalizations of the results in this paragraph, in which the field F of scalars is replaced by a certain kind of ring called a *principal ideal domain*.

6.12 Algebras of Matrices and Linear Operators

Let F be a field. Recall from §1.3 that an *F -algebra* is a structure $(\mathcal{A}, +, \cdot, \star)$ such that $(\mathcal{A}, +, \cdot)$ is an F -vector space, $(\mathcal{A}, +, \star)$ is a ring, and $c \cdot (x \star y) = (c \cdot x) \star y = x \star (c \cdot y)$ for all $x, y \in \mathcal{A}$ and all $c \in F$. (According to this definition, all F -algebras are assumed to be *associative* and have a *multiplicative identity*. There are more general versions of F -algebras, but they will not occur in this book.) A map $g : \mathcal{A} \rightarrow \mathcal{B}$ between two F -algebras is called an *F -algebra homomorphism* iff g is F -linear and a ring homomorphism, i.e.,

$$g(x + y) = g(x) + g(y), \quad g(c \cdot x) = c \cdot g(x), \quad \text{and} \quad g(x \star y) = g(x) \star g(y)$$

for all $x, y \in \mathcal{A}$ and $c \in F$. An F -algebra *isomorphism* is a bijective F -algebra homomorphism.

In linear algebra, there are two premier examples of F -algebras. First, we have the “concrete” F -algebra $M_n(F) = M_{n,n}(F)$ consisting of all square $n \times n$ matrices over the field F . Using the standard matrix operations, $M_n(F)$ is both a vector space (of dimension n^2) and a ring, and $c(AB) = (cA)B = A(cB)$ for all $A, B \in M_n(F)$ and all $c \in F$. The multiplicative identity element is the identity matrix I_n . Second, we have the “abstract” F -algebra $L(V) = L(V, V)$ consisting of all F -linear operators $T : V \rightarrow V$ on an n -dimensional vector space V . We know $L(V)$ is a vector space of dimension n^2 under the standard pointwise operations on linear maps. This vector space becomes a ring (as one readily verifies) if we define the product of $S, T \in L(V)$ to be the composition $S \circ T$. Moreover, id_V is the identity element for $L(V)$.

We have already seen that the *vector spaces* $M_n(F)$ and $L(V)$ are isomorphic (§6.6). Indeed, for each pair of ordered bases X and Y for V , the map $T \mapsto {}_Y[T]_X$ is a vector space isomorphism $L(V) \cong M_n(F)$. Now we make the stronger statement that the *F -algebras* $M_n(F)$ and $L(V)$ are isomorphic. We can obtain *algebra* isomorphisms by demanding that $X = Y$. More specifically, for each ordered basis X of V , consider the vector space isomorphism $M_X = M_{X,X} : L(V) \rightarrow M_n(F)$ given by $M_X(T) = {}_X[T]_X$ for $T \in L(V)$. By (6.6), for all $S, T \in L(V)$,

$$M_X(S \circ T) = {}_X[S \circ T]_X = {}_X[S]_X \ {}_X[T]_X = M_X(S)M_X(T).$$

This says that the bijective linear map M_X preserves products, so it is an algebra isomorphism. From now on, we shall write $[T]_X$ as an abbreviation for ${}_X[T]_X$.

Let us review some previously established facts using the new notation. Suppose $S, T \in L(V)$, $v \in V$, $c \in F$, and $X = (x_1, \dots, x_n)$ is an ordered basis for V . First, by definition,

$$\begin{aligned} A &= [S]_X \text{ iff for all } j \in [n], A^{[j]} = [S(x_j)]_X \text{ iff for all } j \in [n], S(x_j) = L_X(A^{[j]}) \\ &= \sum_{i=1}^n A(i, j)x_i. \end{aligned}$$

Second, M_X preserves the algebraic structure:

$$\begin{aligned} [S + T]_X &= [S]_X + [T]_X; \quad [cS]_X = c[S]_X; \quad [S \circ T]_X = [S]_X [T]_X; \\ [S(v)]_X &= [S]_X [v]_X; \quad [0_V]_X = 0_{n \times n}; \quad [\text{id}_V]_X = I_n. \end{aligned}$$

Third, S is an invertible operator in $L(V)$ iff $[S]_X$ is an invertible matrix in $M_n(F)$, in which case

$$[S^{-1}]_X = ([S]_X)^{-1}.$$

It follows from these remarks that $[S^k]_X = ([S]_X)^k$ holds for all integers $k \geq 0$, and for all negative k too when S is invertible.

An element x in an F -algebra is called *nilpotent* iff $x^m = 0$ for some $m \geq 1$. For example, a square matrix is nilpotent if multiplying the matrix by itself enough times gives the zero matrix. A linear operator on V is nilpotent if applying the operator enough times in succession gives the zero operator on V . From the identity $[S^k]_X = ([S]_X)^k$, we see that $S \in L(V)$ is a nilpotent operator iff $[S]_X \in M_n(F)$ is a nilpotent matrix.

6.13 Similarity of Matrices and Linear Maps

Two $n \times n$ matrices A and A' are called *similar* iff there exists an invertible matrix $P \in M_n(F)$ such that $A' = P^{-1}AP$. Similarity of matrices is readily seen to be an equivalence relation on $M_n(F)$, so $M_n(F)$ is the disjoint union of equivalence classes relative to the similarity relation.

To understand the abstract significance of similarity of matrices, suppose $T \in L(V)$ and X is a fixed ordered basis for V . Let $A = [T]_X$. If we change the basis for V from X to X' , what happens to the representing matrix A for T ? Letting $A' = [T]_{X'}$ and $P = {}_{X'}[\text{id}_V]_{X'}$, (6.6) gives

$$A' = {}_{X'}[\text{id}_V]_{X'}[T]_X{}_{X'}[\text{id}_V]_{X'} = P^{-1}AP.$$

In other words, the matrix $[T]_{X'}$ is similar to the matrix $[T]_X$. Holding X fixed and letting X' vary over all ordered bases for V , we know that P ranges over all invertible $n \times n$ matrices (§6.10). Therefore, the set of all matrices similar to A , namely $\{P^{-1}AP : P \in \text{GL}_n(F)\}$, is precisely the set of all matrices $[T]_{X'}$ that represent the given linear operator T relative to all possible ordered bases X' for V .

Now we consider an “abstract” version of similarity. Two linear maps $T, T' : V \rightarrow V$ are called *similar* iff there exists an invertible linear map $S : V \rightarrow V$ such that $T' = S^{-1} \circ T \circ S$. Similarity of linear maps defines an equivalence relation on the F -algebra $L(V)$. Choosing a fixed ordered basis X for V and applying the isomorphism $M_X : L(V) \cong M_n(F)$, we see immediately that the linear maps T' and T are similar iff the matrices $A' = [T']_X$ and $A = [T]_X$ are similar. As T' ranges over all linear maps similar to T , A' ranges over all matrices similar to A .

To summarize: given a matrix $A = [T]_X \in M_n(F)$, we now have two different interpretations for its equivalence class $\{P^{-1}AP : P \in \text{GL}_n(F)\}$ relative to similarity. First, this equivalence class is the set of all matrices of the form $[T]_{X'}$ as X' ranges over ordered bases for V . Second, this equivalence class also equals the set of all matrices $[T']_X$ where T' ranges over all linear maps similar to T . The first interpretation of equivalence classes is fundamental in linear algebra, for the following reason. If we can find a particularly simple matrix in the similarity class of $A = [T]_X$, then the action of the operator T is correspondingly simple for an appropriately chosen ordered basis X' for V . More specifically, if $A' = P^{-1}AP$ is “nice” (in some sense) when P is the matrix ${}_{X'}[\text{id}]_{X'}$, then the relation $A' = [T]_{X'}$ implies that the action of T on the basis X' will be correspondingly nice. The next section describes some of the “nice” matrices we might hope to find in the equivalence class of A under similarity.

6.14 Diagonalizability and Triangulability

A matrix $A \in M_n(F)$ is called *diagonal* iff $A(i,j) = 0$ for all $i \neq j$; *upper-triangular* iff $A(i,j) = 0$ for all $i > j$; *strictly upper-triangular* iff $A(i,j) = 0$ for all $i \geq j$; *lower-triangular* iff $A(i,j) = 0$ for all $i < j$; and *strictly lower-triangular* iff $A(i,j) = 0$ for all $i \leq j$. For example, the matrices

$$\begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}, \begin{bmatrix} 2 & 1 \\ 0 & 3 \end{bmatrix}, \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, \begin{bmatrix} 2 & 0 \\ 1 & 3 \end{bmatrix}, \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix},$$

are diagonal, upper-triangular, strictly upper-triangular, lower-triangular, and strictly lower-triangular, respectively.

Suppose $T \in L(V)$ where $\dim(V) = n$. The linear operator T is called *diagonalizable* iff there exists an ordered basis $X = (x_1, \dots, x_n)$ for V such that $A = [T]_X$ is a diagonal matrix. In this case,

$$T(x_j) = L_X(A^{[j]}) = \sum_{i=1}^n A(i, j)x_i = A(j, j)x_j \quad \text{for all } j \in [n],$$

so that every $x_j \in X$ is an eigenvector for T with eigenvalue $A(j, j)$. Conversely, if X is an ordered basis such that $T(x_j) = c_j x_j$ for each $x_j \in X$, then $[T]_X$ is the diagonal matrix with c_1, \dots, c_n on the main diagonal. **To summarize:** T is diagonalizable iff there exists an ordered basis X for V consisting of eigenvectors for T .

Next, define $T \in L(V)$ to be *triangulable* (resp. *strictly triangulable*) iff there exists an ordered basis $X = (x_1, \dots, x_n)$ for V such that $A = [T]_X$ is an upper-triangular matrix (resp. strictly upper-triangular matrix). By definition, $T(x_j) = \sum_{i=1}^n A(i, j)x_i$ for all $j \in [n]$. We thereby see that the condition $A(i, j) = 0$ for all $i > j$ is equivalent to the identities $T(x_j) = \sum_{i=1}^j A(i, j)x_i$ for $j \in [n]$. In other words, $[T]_X$ is upper-triangular iff for all $j \in [n]$, $T(x_j)$ is a linear combination of x_1, \dots, x_j . Similarly, $[T]_X$ is strictly upper-triangular iff $T(x_1) = 0$ and each $T(x_j)$ with $2 \leq j \leq n$ is a linear combination of x_1, \dots, x_{j-1} . Analogous considerations show that $[T]_X$ is lower-triangular iff each $T(x_j)$ is a linear combination of x_j, \dots, x_n (similarly for strict lower-triangularity).

We can also phrase these results in terms of similarity of matrices. Given $T \in L(V)$, define $A = [T]_Y$ where Y is any ordered basis for V (e.g., Y might be the standard ordered basis for $V = F^n$). T is diagonalizable iff $A' = [T]_X$ is diagonal for some ordered basis X iff the similarity equivalence class of A contains a diagonal matrix iff $P^{-1}AP$ is diagonal for some invertible matrix $P \in M_n(F)$. Analogously, T is triangulable iff $A' = [T]_X$ is upper-triangular for some ordered basis X iff the similarity equivalence class of A contains an upper-triangular matrix iff $P^{-1}AP$ is upper-triangular for some invertible $P \in M_n(F)$. Accordingly, we say that a matrix A is *diagonalizable* (resp. *triangulable*) iff $P^{-1}AP$ is diagonal (resp. upper-triangular) for some invertible $P \in M_n(F)$.

In Chapters 7 and 8, we will derive theoretical results showing that certain classes of matrices are diagonalizable or triangulable. For now, we merely describe an algorithm for deciding if a specific matrix $A \in M_n(F)$ is diagonalizable and finding a matrix $P \in \mathrm{GL}_n(F)$ such that $P^{-1}AP$ is diagonal. To obtain this algorithm, we consider the left-multiplication map $L^A : F^n \rightarrow F^n$ given by $L^A(v) = Av$ for all column vectors $v \in F^n$. We know every matrix $P^{-1}AP$ has the form $[L^A]_X$ for some ordered basis $X = (x_1, \dots, x_n)$ of V ; moreover, $P = E[\mathrm{id}_V]_X$ where $E = (e_1, \dots, e_n)$ is the standard ordered basis of F^n . The matrix $[L^A]_X$ is diagonal with diagonal entries $c_1, \dots, c_n \in F$ iff $L^A(x_i) = c_i x_i$ for $1 \leq i \leq n$. So, A is diagonalizable iff A has n linearly independent eigenvectors $x_1, \dots, x_n \in F^n$. When this condition holds, we can take P to be the matrix such that $P^{[j]} = [x_j]_E = x_j$ for $1 \leq j \leq n$, so that the columns of P are the linearly independent eigenvectors of A . Then $D = P^{-1}AP$ will be a diagonal matrix with the eigenvalues of A appearing on the diagonal. Note that our “algorithm” for finding D and P assumes that we can compute all the eigenvalues and eigenvectors of A exactly. This may not be true in practice; some numerical methods for approximating the eigenvalues and eigenvectors of A will be covered in Chapter 10.

Next we want to discuss an abstract formulation of triangulability. To do so, we introduce the notion of a flag of subspaces. Given an n -dimensional vector space V , a (*complete*) *flag of subspaces* is a list (V_0, V_1, \dots, V_n) of subspaces of V such that $\{0\} = V_0 \subseteq V_1 \subseteq \dots \subseteq V_n = V$ and $\dim(V_i) = i$ for each i . For example, if $X = (x_1, \dots, x_n)$ is an ordered basis for V and we let V_i be the subspace spanned by x_1, \dots, x_i , then (V_0, \dots, V_n) is a flag of subspaces.

Call this flag the flag of subspaces *associated with* the ordered basis X . Conversely, we can obtain an ordered basis X for V from a flag of subspaces by letting x_i be any element of $V_i \sim V_{i-1}$; note that X is not uniquely determined by the flag. We say that a linear map $T \in L(V)$ *stabilizes* a flag of subspaces iff $T[V_i] \subseteq V_i$ for all i with $0 \leq i \leq n$; T is said to *strictly stabilize* the flag iff $T[V_i] \subseteq V_{i-1}$ for all $i \in [n]$. From our earlier discussion of triangulability, we see that $[T]_X$ is (strictly) upper-triangular iff T (strictly) stabilizes the flag of subspaces associated with X . Conversely, if T (strictly) stabilizes some flag, and we form an ordered basis X from this flag as described above, then $[T]_X$ is (strictly) upper-triangular. **To summarize:** $T \in L(V)$ is (strictly) triangulable iff there exists a flag of subspaces of V (strictly) stabilized by T .

In closing, we remark that a strictly triangulable operator T must be nilpotent. For, letting $V_0 \subseteq V_1 \subseteq \cdots \subseteq V_n = V$ be a flag strictly stabilized by T , and setting $V_i = \{0\}$ for all $i < 0$, we see by induction on $k \geq 1$ that $T^k[V_j] \subseteq V_{j-k}$ for all $j \leq n$. In particular, $T^n[V] = T^n[V_n] \subseteq V_0 = \{0\}$, so that $T^n = 0$ and T is nilpotent. Translating these comments into matrix terms, we see that a strictly upper-triangular matrix A is nilpotent, with $A^n = 0$.

6.15 Block-Triangular Matrices and Invariant Subspaces

We say a matrix $A \in M_n(F)$ is *block-triangular with two blocks* iff there exists $k \in \{1, \dots, n-1\}$ such that for all $i, j \in [n]$ with $i > k$ and $j \leq k$, $A(i, j) = 0$. Such a matrix has the form $\begin{bmatrix} B & C \\ 0 & D \end{bmatrix}$, where $B \in M_k(F)$, $C \in M_{k, n-k}(F)$, $0 \in M_{n-k, k}(F)$, and $D \in M_{n-k}(F)$.

Now let W be a subspace of the vector space V . Given $T \in L(V)$, we say W is a *T -invariant subspace* of V iff $T[W] \subseteq W$, which means that $T(w) \in W$ for all $w \in W$. By linearity, it is equivalent to require that $T(x) \in W$ for all x in some ordered basis for W . If W is T -invariant, then the restricted linear map $T|W : W \rightarrow V$ can be viewed as a linear map with codomain W ; in other words, $T|W$ belongs to the F -algebra $L(W)$. Furthermore, we have an induced map $T' : V/W \rightarrow V/W$ on the quotient vector space V/W , defined by $T'(v + W) = T(v) + W$ for all $v \in V$. This map is well-defined, because for $v, u \in V$, $v + W = u + W$ implies $v - u \in W$, hence $T(v) - T(u) \in W$ by linearity and T -invariance of W , hence $T(v) + W = T(u) + W$. Since T' is evidently linear, T' belongs to the F -algebra $L(V/W)$.

We now relate these abstract considerations to block-triangular matrices. Suppose that W is a T -invariant subspace of V with $\{0\} \neq W \neq V$, and $X = (x_1, \dots, x_n)$ is any ordered basis for V such that $X_1 = (x_1, \dots, x_k)$ is an ordered basis for W . One may check that $X_2 = (x_{k+1} + W, \dots, x_n + W)$ is an ordered basis for V/W (Exercise 51). We will prove that the matrix $A = [T]_X$ is block-triangular with block sizes k and $n-k$; furthermore, the upper-left $k \times k$ block of A is $[T|W]_{X_1}$, whereas the lower-right $(n-k) \times (n-k)$ block of A is $[T']_{X_2}$. Pictorially,

$$A = [T]_X = \begin{bmatrix} [T|W]_{X_1} & C \\ 0 & [T']_{X_2} \end{bmatrix},$$

where C is some $k \times (n-k)$ matrix.

To prove this, let us first compute the column $A^{[j]}$ where $1 \leq j \leq k$. On one hand, since $A = [T]_X$, we have $T(x_j) = \sum_{i=1}^n A(i, j)x_i$, where this expression of $T(x_j)$ in terms of the

basis X is unique. On the other hand, letting $B = [T|W]_{X_1}$, we have

$$T(x_j) = T|W(x_j) = \sum_{i=1}^k B(i,j)x_i = \sum_{i=1}^k B(i,j)x_i + 0x_{k+1} + \cdots + 0x_n.$$

We have now written down two expressions for $T(x_j)$ as linear combinations of the elements in the basis X . Since the expansion in terms of a basis is unique, we conclude that $A(i,j) = B(i,j)$ for $1 \leq i \leq k$, whereas $A(i,j) = 0$ for $k < i \leq n$. These conclusions hold for each $j \in [k]$, so we have confirmed that the upper-left and lower-left blocks of A are given by $B = [T|W]_{X_1}$ and $0 \in M_{n-k,k}(F)$, respectively.

Now consider a column $A^{[j]}$ with $k < j \leq n$. Let $D = [T']_{X_2}$. On one hand, starting from $T(x_j) = \sum_{i=1}^n A(i,j)x_i$, the definition of T' shows that

$$T'(x_j + W) = \sum_{i=1}^n A(i,j)(x_i + W) = \sum_{i=k+1}^n A(i,j)(x_i + W).$$

The last equality holds because $x_i + W = 0 + W$ for $1 \leq i \leq k$, since $x_i \in W$. On the other hand, the definitions of D and X_2 show that

$$T'(x_j + W) = \sum_{i=1}^{n-k} D(i,j-k)(x_{i+k} + W) = \sum_{i=k+1}^n D(i-k,j-k)(x_i + W).$$

Comparing these expressions, we see that $A(i,j) = D(i-k,j-k)$, so that the lower-right block of A is the matrix $D = [T']_{X_2}$.

Conversely, suppose X is any ordered basis for V such that $A = [T]_X$ is block-triangular with diagonal blocks of size k and $n - k$, in this order. Let W be the subspace spanned by $X_1 = (x_1, \dots, x_k)$. For $1 \leq j \leq k$, the block-triangularity of A gives $T(x_j) = \sum_{i=1}^n A(i,j)x_i = \sum_{i=1}^k A(i,j)x_i$, so that $T(x_j) \in W$ for all x_j in X_1 . We conclude that W is a T -invariant subspace of V .

6.16 Block-Diagonal Matrices and Reducing Subspaces

A matrix $A \in M_n(F)$ is *block-diagonal with two blocks* iff there exist $k \in \{1, \dots, n-1\}$, $B \in M_k(F)$, and $D \in M_{n-k}(F)$ with

$$A = \begin{bmatrix} B & 0_{k \times (n-k)} \\ 0_{(n-k) \times k} & D \end{bmatrix}.$$

For $T \in L(V)$, we say that a pair of subspaces W and Z *reduce* T iff W and Z are T -invariant subspaces of V such that $V = W \oplus Z$, i.e., $V = W + Z = \{w+z : w \in W, z \in Z\}$ and $W \cap Z = \{0\}$. Assuming this condition holds, we must have $n = \dim(V) = \dim(W) + \dim(Z)$. Let $X_1 = (x_1, \dots, x_k)$ be an ordered basis for W , and let $X_2 = (x_{k+1}, \dots, x_n)$ be an ordered basis of Z . One readily checks that $X = (x_1, \dots, x_k, x_{k+1}, \dots, x_n)$ is an ordered basis of V . Furthermore, we claim that

$$A = [T]_X = \begin{bmatrix} B & 0 \\ 0 & D \end{bmatrix} \text{ where } B = [T|W]_{X_1} \text{ and } D = [T|Z]_{X_2}.$$

On one hand, if $1 \leq j \leq k$, note x_j lies in the T -invariant subspace W , so $T(x_j) =$

$\sum_{i=1}^n A(i, j)x_i = \sum_{i=1}^k A(i, j)x_i$; also, $T(x_j) = T|W(x_j) = \sum_{i=1}^k B(i, j)x_i$. On the other hand, if $k < j \leq n$, note x_j lies in the T -invariant subspace Z , so $T(x_j) = \sum_{i=1}^n A(i, j)x_i = \sum_{i=k+1}^n A(i, j)x_i$; also, $T(x_j) = T|Z(x_j) = \sum_{i=1}^{n-k} D(i, j-k)x_{i+k}$. The claim follows from these remarks.

Conversely, suppose $A = [T]_X$ is block-diagonal with two diagonal blocks of size k and $n - k$. Let W be the subspace spanned by (x_1, \dots, x_k) , and let Z be the subspace spanned by (x_{k+1}, \dots, x_n) . Since X is an ordered basis, it readily follows that $V = W \oplus Z$. The block-diagonal form of A shows that W and Z are T -invariant subspaces. **To summarize:** there exists a decomposition $V = W \oplus Z$ into two T -invariant subspaces with $\dim(W) = k$ iff there exists an ordered basis X of V such that $[T]_X$ is block-diagonal with blocks of size k and $n - k$.

6.17 Idempotent Matrices and Projections

A matrix $A \in M_n(F)$ is called *idempotent* iff $A^2 = A$. For example, given $0 \leq k \leq n$, the block-diagonal matrix $I_{k,n} = \begin{bmatrix} I_k & 0 \\ 0 & 0 \end{bmatrix} \in M_n(F)$ is idempotent. Analogously, a linear map $T \in L(V)$ is called *idempotent* iff $T \circ T = T$. Applying the F -algebra isomorphism $T \mapsto [T]_X$, where X is any ordered basis of V , we see that a linear map T is idempotent iff $[T]_X$ is an idempotent matrix.

Now suppose V is an F -vector space and W, Z are subspaces of V such that $V = W \oplus Z$. For all $v \in V$, there exist unique $w \in W$ and $z \in Z$ with $v = w + z$. So, we can define a map $P = P_{W,Z} : V \rightarrow V$ by letting $P(v)$ be the unique $w \in W$ appearing in the expression $v = w + z$. We call P the *projection of V onto W along Z* . We claim that P is F -linear and idempotent, with $\text{img}(P) = W$ and $\ker(P) = Z$. Given $v, v' \in V$ and $c \in F$, write $v = w + z$ and $v' = w' + z'$ for unique $w, w' \in W$ and unique $z, z' \in Z$. By definition, $P(v) = w$ and $P(v') = w'$. Now, $v + v' = (w + w') + (z + z')$ with $w + w' \in W$ and $z + z' \in Z$. So the definition of P gives $P(v + v') = w + w' = P(v) + P(v')$. Similarly, $cv = cw + cz$ with $cw \in W$ and $cz \in Z$, so $P(cv) = cw = cP(v)$. Thus P is F -linear. Given $w \in W$, we have $w = w + 0$ with $w \in W$ and $0 \in Z$. So $P(w) = w$ for all $w \in W$. Now, given $v \in V$ with $v = w + z$ as above, we see that $P^2(v) = P(P(v)) = P(w) = w = P(v)$. This holds for all $v \in V$, so $P^2 = P$. Turning to the image of P , the identity $w = P(w)$ for $w \in W$ shows that $W \subseteq \text{img}(P)$. The reverse inclusion $\text{img}(P) \subseteq W$ is immediate from the definition of P , so $\text{img}(P) = W$. As for the kernel, given $z \in Z$, we have $z = 0 + z$ with $0 \in W$ and $z \in Z$. So, $P(z) = 0$, which shows that $Z \subseteq \ker(P)$. Conversely, fix $v \in V$ with $v \in \ker(P)$. Write $v = w + z$ with $w \in W$ and $z \in Z$. Then $w = P(v) = 0$ shows that $v = z \in Z$. So $\ker(P) \subseteq Z$.

Let us compute $[P]_X$, where X is the ordered basis of V obtained by concatenating an ordered basis $X_1 = (x_1, \dots, x_k)$ for W with an ordered basis $X_2 = (x_{k+1}, \dots, x_n)$ for Z . For $1 \leq j \leq k$, $x_j \in W$, so we know that $P(x_j) = x_j = 1x_j + \sum_{i \neq j} 0x_i$. For $k < j \leq n$, $x_j \in Z$, so we know that $P(x_j) = 0 = \sum_{i=1}^n 0x_i$. We now see that $[P]_X = I_{k,n}$. It also follows that W and Z are P -invariant subspaces of V with $P|W = \text{id}_W$ and $P|Z = 0_{L(Z)}$.

Next we show that for every idempotent $T \in L(V)$, there exist unique subspaces W and Z with $V = W \oplus Z$ and $T = P_{W,Z}$. To prove uniqueness, suppose we had subspaces W, Z, W_1, Z_1 with $V = W \oplus Z = W_1 \oplus Z_1$ and $P_{W,Z} = T = P_{W_1,Z_1}$. Then $Z = \ker(P_{W,Z}) = \ker(T) = \ker(P_{W_1,Z_1}) = Z_1$ and $W = \text{img}(P_{W,Z}) = \text{img}(T) = \text{img}(P_{W_1,Z_1}) = W_1$. Now, to prove existence, assume $T \in L(V)$ is idempotent, and define $W = \text{img}(T)$ and $Z = \ker(T)$.

Let us first check that $V = W \oplus Z$. Given $z \in W \cap Z$, we have $z = T(v)$ for some $v \in V$, and also $T(z) = 0$. So $z = T(v) = T(T(v)) = T(z) = 0$, proving $W \cap Z = \{0\}$. Given $v \in V$, note $v = T(v) + (v - T(v))$ where $T(v) \in W = \text{img}(T)$. Since $T(v - T(v)) = T(v) - T(T(v)) = T(v) - T(v) = 0$, we have $v - T(v) \in Z = \ker(T)$. So $v \in W + Z$, and $V = W \oplus Z$. To finish, we check that $T = P_{W,Z}$. Fix $v \in V$, and write $v = w + z$ with $w \in W$ and $z \in Z$. We have $T(z) = 0$ and $w = T(y)$ for some $y \in V$. Applying T to $v = T(y) + z$ gives

$$T(v) = T(T(y) + z) = T(T(y)) + T(z) = T(y) + 0 = T(y) = w = P_{W,Z}(v),$$

as needed. **To summarize:** every idempotent $T \in L(V)$ is the projection $P_{W,Z}$ determined by unique subspaces $W = \text{img}(T)$ and $Z = \ker(T)$, which satisfy $V = W \oplus Z$. T is idempotent iff for some ordered basis X of V and some $k \in \{0, \dots, n\}$, $[T]_X = I_{k,n}$.

6.18 Bilinear Maps and Matrices

Let V be an n -dimensional F -vector space. A map $B : V \times V \rightarrow F$ is called F -bilinear iff for all $v, v', w \in V$ and all $c \in F$, $B(v + v', w) = B(v, w) + B(v', w)$, $B(w, v + v') = B(w, v) + B(w, v')$, and $B(cv, w) = cB(v, w) = B(v, cw)$. It follows from this definition and induction that for all bilinear maps B and all $x_i, y_j \in V$ and $c_i, d_j \in F$,

$$B \left(\sum_{i=1}^r c_i x_i, \sum_{j=1}^s d_j y_j \right) = \sum_{i=1}^r \sum_{j=1}^s c_i d_j B(x_i, y_j). \quad (6.8)$$

Let $\text{BL}(V)$ be the set of all F -bilinear maps on V . One readily checks that the zero map is F -bilinear; the pointwise sum of two F -bilinear maps is F -bilinear; and any scalar multiple of an F -bilinear map is F -bilinear. So, $\text{BL}(V)$ is a subspace of the vector space of all functions from $V \times V$ to F , and is therefore a vector space.

For each ordered basis $X = (x_1, \dots, x_n)$ of V , we will define a map $N_X : \text{BL}(V) \rightarrow M_n(F)$ that will turn out to be a vector space isomorphism. Given $B \in \text{BL}(V)$, we define $N_X(B)$ to be the matrix $[B]_X$ with i, j -entry $B(x_i, x_j)$ for $i, j \in [n]$. We call $[B]_X$ the *matrix of the bilinear map B relative to the ordered basis X* . One checks readily that N_X is an F -linear map. So we need only confirm that N_X is injective and surjective.

To prove injectivity, suppose $B, C \in \text{BL}(V)$ are two bilinear maps with $N_X(B) = N_X(C)$; we must prove $B = C$. Comparing the entries of $[B]_X$ and $[C]_X$, our assumption tells us that $B(x_i, x_j) = C(x_i, x_j)$ for all $i, j \in [n]$. We must show $B(v, w) = C(v, w)$ for all $v, w \in V$. Fix $v, w \in V$, and write $v = \sum_{i=1}^n c_i x_i$ and $w = \sum_{j=1}^n d_j x_j$ for some $c_i, d_j \in F$. Since B and C satisfy (6.8),

$$B(v, w) = \sum_{i=1}^n \sum_{j=1}^n c_i d_j B(x_i, x_j) = \sum_{i=1}^n \sum_{j=1}^n c_i d_j C(x_i, x_j) = C(v, w).$$

To prove surjectivity, let $A \in M_n(F)$ be any fixed matrix; we must find $B \in \text{BL}(V)$ with $N_X(B) = A$. In other words, we must construct an F -bilinear map on V satisfying $B(x_i, x_j) = A(i, j)$ for all $i, j \in [n]$. To do so, we again fix any $v, w \in V$ and write $v = \sum_{i=1}^n c_i x_i$ and $w = \sum_{j=1}^n d_j x_j$ for some $c_i, d_j \in F$. We now set

$$B(v, w) = \sum_{i=1}^n \sum_{j=1}^n c_i d_j A(i, j),$$

which gives a well-defined function since expansions in terms of the ordered basis X are unique. One must now check that B is F -bilinear and $B(x_i, x_j) = A(i, j)$ for $i, j \in [n]$. We prove one identity and leave the others to the reader. Given $v, w \in W$ as above, and also $v' = \sum_{i=1}^n c'_i x_i$ with $c'_i \in F$, note that $v + v' = \sum_{i=1}^n (c_i + c'_i) x_i$. Therefore, the definition of B gives

$$\begin{aligned} B(v + v', w) &= \sum_{i=1}^n \sum_{j=1}^n (c_i + c'_i) d_j A(i, j) \\ &= \sum_{i=1}^n \sum_{j=1}^n c_i d_j A(i, j) + \sum_{i=1}^n \sum_{j=1}^n c'_i d_j A(i, j) = B(v, w) + B(v', w). \end{aligned}$$

To summarize: for each ordered basis X of V , there is an F -vector space isomorphism $\text{BL}(V) \cong M_n(F)$ that sends a bilinear map $B \in \text{BL}(V)$ to the matrix $[B]_X$ with i, j -entry $B(x_i, x_j)$. Therefore, $\dim(\text{BL}(V)) = n^2$, where $n = \dim(V)$.

6.19 Congruence of Matrices

Two matrices $A, A' \in M_n(F)$ are called *congruent* iff there exists an invertible matrix $P \in M_n(F)$ with $A' = P^T AP$. Congruence of matrices defines an equivalence relation on $M_n(F)$, as one readily checks. To understand the abstract significance of congruence of matrices, suppose $B \in \text{BL}(V)$ and we have two ordered bases $X = (x_1, \dots, x_n)$ and $Y = (y_1, \dots, y_n)$ of V . Given that $A = [B]_X$, what happens to A if we change the basis from X to Y ? We will show that A is replaced by a congruent matrix $A' = P^T AP$, where $P = {}_X[\text{id}]_Y$. In particular, as Y ranges over all ordered bases of V , we know that P will range over all invertible matrices in $M_n(F)$. So the set of matrices representing B relative to some ordered basis of V is precisely the set of matrices congruent to A .

Let us now compute $A' = [B]_Y$. We know $y_j = \sum_{i=1}^n P(i, j)x_i$ for all $j \in [n]$. For fixed $r, s \in [n]$, the r, s -entry of A' is

$$\begin{aligned} B(y_r, y_s) &= B\left(\sum_{i=1}^n P(i, r)x_i, \sum_{j=1}^n P(j, s)x_j\right) \\ &= \sum_{i=1}^n \sum_{j=1}^n P(i, r)P(j, s)B(x_i, x_j) \\ &= \sum_{i=1}^n \sum_{j=1}^n P^T(r, i)A(i, j)P(j, s) = (P^T AP)(r, s). \end{aligned}$$

So $A' = P^T AP$ as claimed.

As in the case of similarity, for each congruence class in $M_n(F)$, one would like to find a matrix in that congruence class that is as simple as possible. The answer to this question depends heavily on the field F .

6.20 Real Inner Product Spaces and Orthogonal Matrices

In this section, we consider real vector spaces ($F = \mathbb{R}$). A *real inner product space* consists of a real vector space V and a bilinear map $B \in \text{BL}(V)$ such that $B(v, w) = B(w, v)$ for all $v, w \in V$, and $B(v, v) > 0$ for all nonzero $v \in V$. For $v, w \in V$, we write $\langle v, w \rangle_V = B(v, w)$. For example, $V = \mathbb{R}^n$ is a “concrete” inner product space if we take the bilinear form B to be the *dot product*

$$B(v, w) = \langle v, w \rangle_{\mathbb{R}^n} = v \bullet w = v_1 w_1 + v_2 w_2 + \cdots + v_n w_n \text{ for all } v, w \in \mathbb{R}^n.$$

One sees immediately that $[B]_E = I_n$ where $E = (e_1, \dots, e_n)$ is the standard ordered basis of \mathbb{R}^n . In any inner product space V , the *norm* or *length* of $v \in V$ is defined by $\|v\| = \sqrt{B(v, v)}$. An ordered basis $X = (x_1, \dots, x_n)$ of V is called an *orthonormal basis* iff for all $i, j \in [n]$, $\langle x_i, x_j \rangle_V = \chi(i = j)$, where $\chi(i = j)$ is 1 if $i = j$ and 0 if $i \neq j$. For example, the standard ordered basis E of \mathbb{R}^n is orthonormal.

Given inner product spaces V and W , a linear map $T \in L(V, W)$ is called *orthogonal* iff for all $v, v' \in V$, $\langle T(v), T(v') \rangle_W = \langle v, v' \rangle_V$. To check if a given linear map T is orthogonal, it suffices to verify that $\langle T(x_i), T(x_j) \rangle_W = \langle x_i, x_j \rangle_V$ for all x_i, x_j in a fixed ordered basis X of V . For, supposing this condition holds, let $v = \sum_{i=1}^n c_i x_i$ and $v' = \sum_{j=1}^n d_j x_j$ be any vectors in V , where $c_i, d_j \in \mathbb{R}$. By linearity of T and bilinearity of the inner product, we find that

$$\begin{aligned} \langle T(v), T(v') \rangle_W &= \left\langle T\left(\sum_{i=1}^n c_i x_i\right), T\left(\sum_{j=1}^n d_j x_j\right) \right\rangle_W = \left\langle \sum_{i=1}^n c_i T(x_i), \sum_{j=1}^n d_j T(x_j) \right\rangle_W \\ &= \sum_{i=1}^n \sum_{j=1}^n c_i d_j \langle T(x_i), T(x_j) \rangle_W = \sum_{i=1}^n \sum_{j=1}^n c_i d_j \langle x_i, x_j \rangle_V \\ &= \left\langle \sum_{i=1}^n c_i x_i, \sum_{j=1}^n d_j x_j \right\rangle_V = \langle v, v' \rangle_V. \end{aligned}$$

Let $O(V)$ denote the set of all orthogonal linear maps from V to V . $O(V)$ is not a *subspace* of the vector space $L(V)$, but $O(V)$ is a *subgroup* of the group $\text{GL}(V)$ of all invertible linear maps $T \in L(V)$; $O(V)$ is called the *orthogonal group on V* . We check that $O(V) \subseteq \text{GL}(V)$ and let the reader confirm the three closure conditions in the definition of a subgroup (see §1.4). Given $T \in O(V)$, we know T is a linear map from V to V . To see that T is invertible, it therefore suffices to check that $\ker(T) = \{0_V\}$. Given $v \in \ker(T)$, orthogonality of T gives $\langle v, v \rangle_V = \langle T(v), T(v) \rangle_V = \langle 0, 0 \rangle_V = 0$, which forces $v = 0$ by the definition of an inner product space.

We have seen that each ordered basis $X = (x_1, \dots, x_n)$ of V induces a vector space isomorphism $L_X : \mathbb{R}^n \rightarrow V$. We now ask when L_X will be an orthogonal map (which can be regarded as an isomorphism of inner product spaces). As remarked above, L_X is orthogonal iff $\langle L_X(e_i), L_X(e_j) \rangle_V = \langle e_i, e_j \rangle_{\mathbb{R}^n}$ for all $i, j \in [n]$ iff $\langle x_i, x_j \rangle_V = \chi(i = j)$ for all $i, j \in [n]$. So, L_X is an inner product space isomorphism $\mathbb{R}^n \cong V$ iff X is an orthonormal basis of V .

Let $X = (x_1, \dots, x_n)$ be an orthonormal basis for the inner product space V . What is the image of the orthogonal group $O(V)$ under the \mathbb{R} -algebra isomorphism from $L(V)$ to $M_n(\mathbb{R})$ that sends $T \in L(V)$ to $[T]_X \in M_n(\mathbb{R})$? To answer this question, fix $T \in L(V)$ with matrix $A = [T]_X$. Note T is in $O(V)$ iff $\langle T(x_i), T(x_j) \rangle_V = \langle x_i, x_j \rangle_V$ for all $i, j \in [n]$ iff $\langle L_X(A^{[i]}), L_X(A^{[j]}) \rangle_V = \chi(i = j)$ for all $i, j \in [n]$ iff $\langle A^{[i]}, A^{[j]} \rangle_{\mathbb{R}^n} = \chi(i = j)$ for all $i, j \in [n]$. Let us call a matrix $A \in M_n(\mathbb{R})$ *orthogonal* iff A satisfies the condition just

stated, and let $O_n(\mathbb{R})$ be the set of all orthogonal matrices in $M_n(\mathbb{R})$. We have just shown that *the linear map T is orthogonal iff its matrix A (relative to an orthonormal basis) is orthogonal.*

Our definition of an orthogonal matrix A states that the column vectors $A^{[1]}, \dots, A^{[n]}$ must be orthonormal. Geometrically, this condition says that each column of A is a unit vector (has norm 1), and any two distinct columns of A are perpendicular. We can restate this condition by noting that $\langle A^{[i]}, A^{[j]} \rangle_{\mathbb{R}^n} = A^{[i]} \bullet A^{[j]}$ is the i, j -entry of the matrix $A^T A$. So $A^{[i]} \bullet A^{[j]} = \chi(i = j)$ for all $i, j \in [n]$ iff $(A^T A)(i, j) = \chi(i = j)$ for all $i, j \in [n]$ iff $A^T A = I_n$ iff A is invertible with $A^{-1} = A^T$ (see §4.6) iff $AA^T = I_n$ iff $(AA^T)(i, j) = \chi(i = j)$ for all $i, j \in [n]$ iff $A_{[i]} \bullet A_{[j]} = \chi(i = j)$ for all $i, j \in [n]$. It now follows that the column vectors $A^{[1]}, \dots, A^{[n]}$ are orthonormal iff A is orthogonal iff A^T is orthogonal iff the row vectors $A_{[1]}, \dots, A_{[n]}$ are orthonormal.

Consider a transition matrix $P = {}_Y[\text{id}]_X$ between two *orthonormal* bases of the inner product space V . By orthonormality of Y , the map L_Y^{-1} sending $v \in V$ to $[v]_Y$ is orthogonal. By orthonormality of X , $\langle x_i, x_j \rangle_V = \chi(i = j)$ for $i, j \in [n]$. It follows that the columns $[x_1]_Y, \dots, [x_n]_Y$ of P are orthonormal in \mathbb{R}^n , so that P is an *orthogonal* matrix. One can check that every orthogonal matrix in $O_n(\mathbb{R})$ arises in this way for some choice of orthonormal bases X and Y . Now, define two matrices $A, A' \in M_n(\mathbb{R})$ to be *orthogonally similar* iff $A' = P^{-1}AP = P^TAP$ for some orthogonal $P \in O_n(\mathbb{R})$. Orthogonal similarity is an equivalence relation on $M_n(\mathbb{R})$. Given $T \in L(V)$, the equivalence class of $A = [T]_X$ consists of all possible matrices that represent T relative to different choices of *orthonormal* bases for V .

6.21 Complex Inner Product Spaces and Unitary Matrices

Now we consider the complex version of inner product spaces ($F = \mathbb{C}$). A *complex inner product space* consists of a complex vector space V and a map $B : V \times V \rightarrow \mathbb{C}$, denoted $B(v, w) = \langle v, w \rangle_V$, satisfying the following identities for all $v, v', w \in V$ and all $c \in \mathbb{C}$:

$$\begin{aligned} \langle v + v', w \rangle_V &= \langle v, w \rangle_V + \langle v', w \rangle_V; & \langle cv, w \rangle_V &= c \langle v, w \rangle_V; & \langle w, v \rangle_V &= \overline{\langle v, w \rangle_V}; \\ \langle v, v \rangle &\in \mathbb{R}^+ \text{ for } v \neq 0. \end{aligned}$$

The notation \bar{z} denotes the complex conjugate of the complex number z , i.e., $\overline{a + ib} = a - ib$ for $a, b \in \mathbb{R}$. It follows from the preceding identities that

$$\langle w, v + v' \rangle_V = \langle w, v \rangle_V + \langle w, v' \rangle_V \text{ and } \langle w, cv \rangle_V = \bar{c} \langle w, v \rangle_V.$$

As before, an ordered basis $X = (x_1, \dots, x_n)$ of V is called *orthonormal* iff $\langle x_i, x_j \rangle_V = \chi(i = j)$ for all $i, j \in [n]$. For example, \mathbb{C}^n is a complex inner product space with inner product

$$\langle v, w \rangle_{\mathbb{C}^n} = v_1 \overline{w_1} + v_2 \overline{w_2} + \cdots + v_n \overline{w_n} \text{ for all } v, w \in \mathbb{C}^n,$$

and the standard ordered basis $E = (e_1, \dots, e_n)$ of \mathbb{C}^n is orthonormal.

For complex inner product spaces V and W , $T \in L(V, W)$ is called a *unitary map* iff $\langle T(v), T(v') \rangle_W = \langle v, v' \rangle_V$ for all $v, v' \in V$. As in the real case, it suffices to check that $\langle T(x_i), T(x_j) \rangle_W = \langle x_i, x_j \rangle_V$ for all x_i, x_j in an ordered basis X of V . Let $U(V)$ be the set of all unitary maps in $L(V)$; $U(V)$ is a subgroup of $GL(V)$.

One checks as in the real case that the isomorphism $L_X : \mathbb{C}^n \rightarrow V$ is a unitary map iff X

is an orthonormal basis of V . Furthermore, given an orthonormal basis X and a linear map $T \in L(V)$, T lies in $U(V)$ iff $A = [T]_X$ is a matrix with orthonormal columns in \mathbb{C}^n . Call a matrix with this property *unitary*, and let $U_n(\mathbb{C})$ be the set of all such matrices in $M_n(\mathbb{C})$. Recall that the conjugate-transpose of A is the matrix A^* with i, j -entry $A^*(i, j) = \overline{A(j, i)}$. The orthonormality of the columns of A in \mathbb{C}^n is equivalent to the matrix identity $A^*A = I_n$, which is equivalent to $AA^* = I_n$ (see §4.6). So, A is unitary iff A is invertible with $A^{-1} = A^*$ iff the rows of A are orthonormal in \mathbb{C}^n .

As in the real case, a transition matrix $P = {}_Y[\text{id}]_X$ between two orthonormal bases of a complex inner product space V must be unitary. Moreover, all unitary matrices have this form for some choice of orthonormal bases X and Y . Define two matrices $A, A' \in M_n(\mathbb{C})$ to be *unitarily similar* iff $A' = P^{-1}AP = P^*AP$ for some unitary $P \in U_n(\mathbb{C})$. Unitary similarity is an equivalence relation on $M_n(\mathbb{C})$. Given $T \in L(V)$, the equivalence class of $A = [T]_X$ consists of all possible matrices that represent T relative to different choices of *orthonormal* bases for V .

6.22 Summary

1. *Notation.* Table 6.2 recalls the notation used in this chapter for various vector spaces, algebras, and groups.

TABLE 6.2

Summary of Notation.

Symbol	Meaning
F^n	vector space of n -tuples (column vectors) with entries in F
$M_{m,n}(F)$	vector space of $m \times n$ matrices with entries in F
$M_n(F)$	F -algebra of $n \times n$ matrices
$GL_n(F)$	group of invertible $n \times n$ matrices
$O_n(\mathbb{R})$	group of orthogonal $n \times n$ real matrices ($A^{-1} = A^T$)
$U_n(\mathbb{C})$	group of unitary $n \times n$ complex matrices ($A^{-1} = A^*$)
V, W	abstract F -vector spaces
$L(V, W)$	vector space of F -linear maps $T : V \rightarrow W$
$L(V)$	F -algebra of F -linear maps $T : V \rightarrow V$
$BL(V)$	vector space of bilinear maps $B : V \times V \rightarrow F$
$GL(V)$	group of invertible F -linear maps $T : V \rightarrow V$
$O(V)$	group of orthogonal linear maps on real inner product space V
$U(V)$	group of unitary linear maps on complex inner product space V

2. *Main Isomorphisms.* Table 6.3 reviews the isomorphisms between abstract algebraic structures and concrete versions of these structures covered in this chapter.
3. *Linear Combination Maps.* Given any list $X = (x_1, \dots, x_n)$ of vectors in some abstract F -vector space V , we have a linear map $L_X : F^n \rightarrow V$ given by $L_X(c_1, \dots, c_n) = \sum_{i=1}^n c_i x_i$ for $c_i \in F$. L_X is surjective iff X spans V , L_X is injective iff X is linearly independent, and L_X is bijective iff X is an ordered basis for V . An F -vector space V has dimension n iff there exists a vector space isomorphism $V \cong F^n$; any two n -dimensional F -vector spaces are isomorphic. The function $X \mapsto L_X$ is a bijection from the set of ordered bases for V onto the

TABLE 6.3

Main Isomorphisms between Abstract Linear-Algebraic Objects and Concrete Matrix-Theoretic Objects.

Isomorphism	Type of Iso.	Isomorphism Depends On:
$V \cong F^n$	vector space	ordered basis X of V
$L(V, W) \cong M_{m,n}(F)$	vector space	ordered bases X of V and Y of W
$L(V) \cong M_n(F)$	F -algebra	ordered basis X of V
$BL(V) \cong M_n(F)$	vector space	ordered basis X of V
$GL(V) \cong GL_n(F)$	group	ordered basis X of V
$V \cong \mathbb{R}^n$ or \mathbb{C}^n	inner prod. space	orthonormal basis X of V
$O(V) \cong O_n(\mathbb{R})$	group	orthonormal basis X of V
$U(V) \cong U_n(\mathbb{C})$	group	orthonormal basis X of V

set of vector space isomorphisms $L : F^n \cong V$. The inverse bijection sends such an isomorphism L to the ordered basis $(L(e_1), \dots, L(e_n))$, where e_1, \dots, e_n are the standard basis vectors in F^n .

4. *Vector Spaces of n -tuples vs. Abstract Vectors.* If $X = (x_1, \dots, x_n)$ is an ordered basis for an abstract F -vector space V , the coordinate map $v \mapsto [v]_X$ is an F -vector space isomorphism $V \cong F^n$. In detail, for all $c, c_i \in F$ and all $v, w \in V$:

$$(c_1, \dots, c_n) = [v]_X \Leftrightarrow v = L_X(c_1, \dots, c_n) \Leftrightarrow v = c_1x_1 + \dots + c_nx_n;$$

$$[v + w]_X = [v]_X + [w]_X; \quad [cv]_X = c[v]_X.$$

5. *Vector Spaces of Matrices and Linear Maps.* If $X = (x_1, \dots, x_n)$ and $Y = (y_1, \dots, y_m)$ are ordered bases for F -vector spaces V and W , then the map $T \mapsto {}_Y[T]_X$ is an F -vector space isomorphism $L(V, W) \cong M_{m,n}(F)$. In detail, for all $A \in M_{m,n}(F)$, $S, T \in L(V, W)$, $c \in F$:

$$A = {}_Y[T]_X \Leftrightarrow \forall j, A^{[j]} = [T(x_j)]_Y \Leftrightarrow \forall j, T(x_j) = L_Y(A^{[j]}) = \sum_{i=1}^n A(i, j)y_i;$$

$${}_Y[S + T]_X = {}_Y[S]_X + {}_Y[T]_X; \quad {}_Y[cT]_X = c({}_Y[T]_X).$$

Moreover, if $R \in L(W, U)$ where U has ordered basis Z , then ${}_Z[R \circ T]_X = {}_Z[R]_Y {}_Y[T]_X$.

6. *Algebras of Matrices and Operators.* For each ordered basis $X = (x_1, \dots, x_n)$ of an F -vector space V , the map $T \mapsto [T]_X$ is an F -algebra isomorphism $L(V) \cong M_n(F)$. In detail, for all $A \in M_n(F)$, $S, T \in L(V)$, $c \in F$, $v \in V$:

$$A = [T]_X \Leftrightarrow \forall j, A^{[j]} = [T(x_j)]_X \Leftrightarrow \forall j, T(x_j) = L_X(A^{[j]}) = \sum_{i=1}^n A(i, j)x_i;$$

$$[S + T]_X = [S]_X + [T]_X; \quad [cS]_X = c[S]_X; \quad [S \circ T]_X = [S]_X [T]_X;$$

$$[S(v)]_X = [S]_X [v]_X; \quad [0_{L(V)}]_X = 0_{M_n(F)}; \quad [\text{id}_V]_X = I_n;$$

$$\forall k \in \mathbb{N}, [S^k]_X = ([S]_X)^k; \text{ this holds for all } k \in \mathbb{Z} \text{ when } S \text{ is invertible.}$$

7. *Transition Matrices.* For ordered bases X and Y of V , the *transition matrix from X to Y* is the matrix $A = {}_Y[\text{id}_V]_X$ defined by $A^{[j]} = [x_j]_Y$ for $j \in [n]$. We have $x_j = \sum_{i=1}^n A(i, j)y_i$, $A^{-1} = {}_X[\text{id}]_Y$, and ${}_X[\text{id}]_X = I_n$. For each fixed X , the map $Y \mapsto {}_Y[\text{id}]_X$ is a bijection from the set of ordered bases of V to the set of invertible matrices $\text{GL}_n(F)$. For each fixed Y , the map $X \mapsto {}_Y[\text{id}]_X$ is also a bijection between these sets. If $P = {}_Y[\text{id}]_X$ and $T \in L(V)$, then $[T]_X = P^{-1}[T]_Y P$.
8. *Bilinear Maps and Congruence of Matrices.* For each ordered basis X of V , there is a vector space isomorphism $\text{BL}(V) \cong M_n(F)$ that sends a bilinear map B to the matrix $[B]_X$ with i, j -entry $B(x_i, x_j)$. Changing the basis from X to Y replaces $A = [B]_X$ by $A' = [B]_Y = P^T AP$, where $P = {}_X[\text{id}]_Y$. So, the set of all matrices congruent to A is the set of matrices representing B relative to some ordered basis of V .
9. *Inner Product Spaces.* For each orthonormal basis X of a real (resp. complex) inner product space V , the map $v \mapsto [v]_X$ (and its inverse L_X) are isomorphisms preserving the inner product. Orthogonal (resp. unitary) maps correspond to orthogonal (resp. unitary) matrices under the isomorphism $T \mapsto [T]_X$, for $T \in L(V)$. The transition matrix between two orthonormal bases is orthogonal (resp. unitary). The equivalence class of a matrix $A = [T]_X$ under orthogonal (resp. unitary) similarity is the set of all matrices that represent T as we vary the orthonormal basis X .
10. *Equivalence Relations on Matrices.* Table 6.4 summarizes some commonly used equivalence relations on rectangular and square matrices.

TABLE 6.4

Equivalence Relations on Matrices.

Eq. Relation	Definition of $A \sim A'$	Abstract Significance
col. equivalence	$\exists P \in \text{GL}_n(F), A' = AP$	changes input basis in ${}_Y[T]_X$
row equivalence	$\exists Q \in \text{GL}_m(F), A' = QA$	changes output basis in ${}_Y[T]_X$
row/col equiv.	$\exists P \in \text{GL}_n(F), \exists Q \in \text{GL}_m(F), A' = QAP$	changes input/output bases
similarity	$\exists P \in \text{GL}_n(F), A' = P^{-1}AP$	go from $[T]_X$ to $[T]_Y$
congruence	$\exists P \in \text{GL}_n(F), A' = P^T AP$	change basis for bilinear map
orthogonal sim.	$\exists P \in O_n(\mathbb{R}), A' = P^{-1}AP = P^T AP$	orthonormal basis change ($F = \mathbb{R}$)
unitary sim.	$\exists P \in U_n(\mathbb{C}), A' = P^{-1}AP = P^* AP$	orthonormal basis change ($F = \mathbb{C}$)

11. *Diagonalizability and Triangulability.* A linear map $T \in L(V)$ is diagonalizable iff $[T]_X$ is diagonal for some ordered basis X of V iff V has an ordered basis consisting of eigenvectors of T iff $A = [T]_Y$ is similar to a diagonal matrix for all ordered bases Y of V iff $P^{-1}AP$ is diagonal for some invertible $P \in M_n(F)$. A *flag of subspaces* of V is a chain of subspaces $\{0\} = V_0 \subseteq V_1 \subseteq \dots \subseteq V_n = V$ with $\dim(V_i) = i$ for all i . T is triangulable iff $[T]_X$ is upper-triangular for some ordered basis X of V iff T stabilizes some flag of subspaces of V iff $A = [T]_Y$ is similar to an upper-triangular matrix for all ordered bases Y of V iff $P^{-1}AP$ is upper-triangular for some invertible $P \in M_n(F)$.
12. *Reducibility and Invariant Subspaces.* Given $T \in L(V)$ and an ordered basis $X = (x_1, \dots, x_n)$ of V , the matrix $A = [T]_X$ has the block-triangular form

$\begin{bmatrix} B_{k \times k} & C \\ 0 & D_{(n-k) \times (n-k)} \end{bmatrix}$ iff $X_1 = (x_1, \dots, x_k)$ spans a T -invariant subspace W of V (which is a subspace with $T(w) \in W$ for all $w \in W$). In this case, $B = [T|W]_{X_1}$ and $D = [T']_{X'_2}$, where $T|W \in L(W)$ is the restriction of T to W , $T' : V/W \rightarrow V/W$ in $L(V/W)$ is given by $T'(v + W) = T(v) + W$ for $v \in V$, and $X'_2 = (x_{k+1} + W, \dots, x_n + W)$. Moreover, $A = [T]_X$ is block-diagonal with two diagonal blocks of size k and $n - k$ iff $C = 0$ iff $X_1 = (x_1, \dots, x_k)$ and $X_2 = (x_{k+1}, \dots, x_n)$ both span T -invariant subspaces of V iff T is reduced by the pair of subspaces W and Z spanned by X_1 and X_2 . In this case, $B = [T|W]_{X_1}$ and $D = [T|Z]_{X_2}$.

13. *Idempotent Matrices and Projections.* A linear map $T \in L(V)$ is *idempotent* iff $T \circ T = T$. Given $V = W \oplus Z$, the projection on W along Z is the map $P_{W,Z}$ that sends each $v \in V$ to the unique $w \in W$ such that $v = w + z$ for some $z \in Z$. The map $P_{W,Z}$ is idempotent with kernel Z and image W , and $[P_{W,Z}]_X = I_{k,n}$, where X is obtained by concatenating ordered bases for W and Z . Every idempotent $T \in L(V)$ has the form $P_{W,Z}$ for unique subspaces W, Z with $V = W \oplus Z$; in fact, $W = \text{img}(T)$ and $Z = \ker(T)$.
14. *Computational Procedures.* Assume $X = (x_1, \dots, x_n)$ and $Y = (y_1, \dots, y_m)$ are ordered bases of F -vector spaces V and W .

- **To compute coordinates** $[v]_X$ when given $v \in V$ and X : solve the equation $v = c_1x_1 + \dots + c_nx_n$ for the unknown c_j 's in F ; and output the answer $[v]_X = (c_1, \dots, c_n)$.
- **To compute the matrix** ${}_Y[T]_X$ **of a linear map** when given T and X and Y : for each $j \in [n]$, find $T(x_j)$; compute the coordinates of this vector relative to Y , namely $[T(x_j)]_Y$; and write down these coordinates as the j 'th column of the matrix.
- **To find** $[T]_X$ given T and X : perform the algorithm in the preceding item taking $Y = X$.
- **To find the transition matrix from X to Y :** compute the coordinates $[x_j]_Y$ of each x_j relative to Y and write them down as the j 'th column of the transition matrix. Alternatively, compute the matrix inverse of the transition matrix from Y to X if the latter matrix is already available. (For instance, this occurs when X is the standard ordered basis for F^n and each $y_j \in Y$ is presented as a column vector, so $[y_j]_X = y_j$ for all j .)
- **To change coordinates** from $[v]_X$ to $[v]_Y$ where $[v]_X$ and X and Y are given: find the transition matrix ${}_Y[\text{id}]_X$ from X to Y ; compute the matrix-vector product $[v]_Y = {}_Y[\text{id}]_X [v]_X$.
- **To find the matrix of a linear map relative to a new basis** where $[T]_X$ and X and Y are given: find the transition matrix $P = {}_X[\text{id}]_Y$ and its inverse $P^{-1} = {}_Y[\text{id}]_X$; compute the matrix product $[T]_Y = P^{-1}[T]_X P$.
- **To diagonalize a diagonalizable linear operator** T : compute an ordered basis X for V consisting of eigenvectors for T (e.g., by finding roots of the characteristic polynomial and solving linear equations); output $[T]_X$, which is the diagonal matrix with the eigenvalues corresponding to each $x_j \in X$ appearing on the main diagonal.
- **To diagonalize a diagonalizable matrix** A : use the preceding item to diagonalize the left multiplication map $L^A : F^n \rightarrow F^n$. If E is the standard basis for F^n and $P = {}_E[\text{id}]_X$ is the matrix whose columns are the

eigenvectors of A , then $P^{-1}AP$ is a diagonal matrix with the eigenvalues of A on the main diagonal.

- **To compute the matrix $[B]_X$ of a bilinear map B :** let the i, j -entry be $B(x_i, x_j)$ for all $i, j \in [n]$.
- **To find the matrix of a bilinear map relative to a new basis:** compute $P = {}_X[\text{id}]_Y$ and $[B]_Y = P^T[B]_X P$.

Our dictionary for translating between matrix theory and linear algebra is presented in Tables 6.2, 6.3, 6.4, and 6.5.

TABLE 6.5

Dictionary Connecting Concrete Matrix Theory and Abstract Linear Algebra.

Concept in Matrix Theory	Connecting Formula	Concept in Linear Algebra
concrete vector space F^n of n -tuples	isomorphisms $L_X : F^n \cong V$ $L_X(c_1, \dots, c_n) = \sum_{i=1}^n c_i x_i$	abstract F -vector space V with ordered basis $X = (x_1, \dots, x_n)$
n -tuple $(c_1, \dots, c_n) \in F^n$	$[v]_X = (c_1, \dots, c_n)$ iff $v = c_1 x_1 + \dots + c_n x_n$	abstract vector $v \in V$
componentwise operations on n -tuples	$[v+w]_X = [v]_X + [w]_X$, $[cv]_X = c[v]_X$	abstract vector addition and scalar multiplication in V
set of all isomorphisms $L : F^n \cong V$	bijections: $X \mapsto L_X$ and $L \mapsto (L(e_1), \dots, L(e_n))$	set of all ordered bases X of V
vector space $M_{m,n}(F)$ of $m \times n$ matrices	isomorphisms $T \mapsto {}_Y[T]_X$ indexed by ordered bases	vector space $L(V, W)$ of linear maps from V to W
matrix $A \in M_{m,n}(F)$ with columns $A^{[j]}$ and entries $A(i, j)$	$A = {}_Y[T]_X$ iff $A^{[j]} = [T(x_j)]_Y$ iff $\forall j, T(x_j) = \sum_{i=1}^m A(i, j) y_i$	linear map $T \in L(V, W)$ with X an ordered basis for V and Y an ordered basis for W
matrix addition $(A + B)(i, j) = A(i, j) + B(i, j)$	${}_Y[S + T]_X = {}_Y[S]_X + {}_Y[T]_X$	addition of linear maps $(S + T)(v) = S(v) + T(v)$
scalar multiple of a matrix $(cA)(i, j) = c(A(i, j))$	${}_Y[cT]_X = c({}_Y[T]_X)$	scalar multiple of linear map $(cT)(v) = c \cdot T(v)$
matrix/vector multiplication $(Av)_i = \sum_{j=1}^n A(i, j)v_j$	$[T(v)]_Y = {}_Y[T]_X [v]_X$	evaluation of linear map T applied to v gives $T(v)$
matrix multiplication $(BA)(i, j) = \sum_k B(i, k)A(k, j)$	$z[S \circ T]_X = z[S]_Y {}_Y[T]_X$	composition of linear maps $(S \circ T)(v) = S(T(v))$
inverse of a matrix $AA^{-1} = I_n = A^{-1}A$	$x[T^{-1}]_Y = ({}_Y[T]_X)^{-1}$	inverse of a map $T \circ T^{-1} = \text{id}_W$, $T^{-1} \circ T = \text{id}_V$
matrix transpose $A^T(i, j) = A(j, i)$	$x^*[T^*]_{Y^*} = ({}_Y[T]_X)^T$ (X^* , Y^* are dual bases)	dual map $T^* \in L(W^*, V^*)$ $T^*(g) = g \circ T : V \rightarrow F$
$m \times n$ matrix A	$E^*[L^A]_E = A$ (E, E' are standard bases)	left multiplication by A $L^A : F^n \rightarrow F^m$, $L^A(v) = Av$
F -algebra $M_n(F)$	isomorphisms $T \mapsto [T]_X$	F -algebra $L(V)$
identity matrix I_n	$[\text{id}_V]_X = I_n$	identity map $\text{id}_V : V \rightarrow V$ $\text{id}_V(v) = v$ for $v \in V$
zero matrix 0_n	$[0_V]_X = 0_n$	zero map $0_V : V \rightarrow V$ $0_V(v) = 0 \in V$ for $v \in V$

6.23 Exercises

Unless otherwise stated, assume that F is a field, V is an F -vector space with ordered basis $X = (x_1, \dots, x_n)$, and W is an F -vector space with ordered basis $Y = (y_1, \dots, y_m)$. Exercises involving \mathbb{R}^3 , $P_{\leq 3}$, $M_2(\mathbb{R})$, and \mathbb{C} will use the notation from Table 6.1.

1. (a) For each vector space V in Table 6.1, prove that X_1 is an ordered basis for V and thereby compute $\dim(V)$. (b) For each V in the table, prove that X_2 and X_3 are ordered bases for V . (By (a), it suffices to prove that each list is linearly independent.)
2. For each vector $v \in \mathbb{R}^3$, compute $[v]_{X_1}$, $[v]_{X_2}$, and $[v]_{X_3}$.
 - (a) $v = (1, 2, 4)$; (b) $v = (-2, 5, 2)$; (c) $v = (0, 0, c)$.
3. For each vector $v \in P_{\leq 3}$, compute $[v]_{X_1}$, $[v]_{X_2}$, and $[v]_{X_3}$.
 - (a) $v = 1 + t + t^2 + t^3$; (b) $v = (t+1)^2$; (c) $v = t(3t-4)(2t+1)$.
4. For each vector $v \in M_2(\mathbb{R})$, compute $[v]_{X_1}$, $[v]_{X_2}$, and $[v]_{X_3}$.
 - (a) $v = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$; (b) $v = \begin{bmatrix} 4 & 4 \\ -1 & -1 \end{bmatrix}$; (c) $v = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$.
5. For each vector $v \in \mathbb{C}$, compute $[v]_{X_1}$, $[v]_{X_2}$, and $[v]_{X_3}$.
 - (a) $v = -i$; (b) $v = 7 - 5i$; (c) $v = e^{\pi i/4}$.
6. (a) Which vector $v \in \mathbb{R}^3$ has $[v]_{X_2} = (3, 2, 1)$? (b) Which vector $v \in P_{\leq 3}$ has $[v]_{X_3} = (2, -1, 0, 5)$? (c) Which vector $v \in \mathbb{C}$ has $[v]_{X_3} = (2, -3)$?
7. (a) For all $c \in \mathbb{R}$, find $[(t+c)^3]_{X_1}$. (b) For all $c \in \mathbb{R}$, find $[(t+c)^3]_{X_2}$.
8. (a) Find a nonzero $A \in M_2(\mathbb{R})$ such that $[A]_{X_1} = 6[A]_{X_3}$. (b) Find all $c \in \mathbb{R}$ such that for some nonzero $A \in M_2(\mathbb{R})$, $[A]_{X_1} = c[A]_{X_3}$.
9. Let X' be obtained from X by switching the positions of x_i and x_j . (a) Show X' is an ordered basis of V . (b) For $v \in V$, how is $[v]_{X'}$ related to $[v]_X$?
10. Let X' be obtained from X by replacing x_i by bx_i , where $b \in F$ is nonzero. (a) Show X' is an ordered basis of V . (b) For $v \in V$, how is $[v]_{X'}$ related to $[v]_X$?
11. Let X' be obtained from X by replacing x_i by $x_i + ax_j$, where $a \in F$ and $j \neq i$. (a) Show X' is an ordered basis of V . (b) For $v \in V$, how is $[v]_{X'}$ related to $[v]_X$?
12. Define $T : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ by $T(a, b, c) = (a + b - c, 3b + 2c)$ for $a, b, c \in \mathbb{R}$. Let $Y = ((1, 0), (0, 1))$ and $Z = ((1, 3), (-1, 2))$. (a) Find ${}_Y[T]_{X_1}$ and ${}_Z[T]_{X_1}$. (b) Find ${}_Y[T]_{X_2}$ and ${}_Z[T]_{X_2}$. (c) There is an ordered basis W of \mathbb{R}^2 such that ${}_W[T]_{X_3} = \begin{bmatrix} 3 & 1 & b \\ -1 & 1 & c \end{bmatrix}$. Find W , b , and c .
13. Define $T : P_{\leq 3} \rightarrow \mathbb{R}^3$ by $T(f) = (f(2), f'(2), \int_0^2 f(t) dt)$ for $f \in P_{\leq 3}$. (a) Check that T is \mathbb{R} -linear. (b) For $Y = (e_1, e_2, e_3)$, compute ${}_Y[T]_{X_1}$, ${}_Y[T]_{X_2}$, and ${}_Y[T]_{X_3}$. (c) Repeat (b) for $Y = ((1, 1, 1), (1, 2, 4), (1, 3, 9))$.
14. Let $\text{tr} : M_2(\mathbb{R}) \rightarrow \mathbb{R}$ be the trace map, defined by $\text{tr}(A) = A(1, 1) + A(2, 2)$ for $A \in M_2(\mathbb{R})$. Let $Z = (1)$. (a) Explain why $\text{tr} \in M_2(\mathbb{R})^*$. (b) Find ${}_Z[\text{tr}]_{X_1}$, ${}_Z[\text{tr}]_{X_2}$, and ${}_Z[\text{tr}]_{X_3}$. (c) Let $T : M_2(\mathbb{R}) \rightarrow M_2(\mathbb{R})$ be the transpose map. Show that $T^*(\text{tr}) = \text{tr}$.
15. Define $T : \mathbb{C} \rightarrow \mathbb{C}$ by $T(z) = iz$ for $z \in \mathbb{C}$. (a) Compute $[T]_{X_1}$, $[T]_{X_2}$, and $[T]_{X_3}$. (b) Square each matrix found in (a), and discuss the answers. (c) Find ${}_{X_3}[T]_{X_1}$ and ${}_{X_1}[T]_{X_3}$.

16. Let $A = \begin{bmatrix} 2 & 1 \\ -1 & 3 \end{bmatrix}$. Define three maps $\lambda_A, \rho_A, \kappa_A : M_2(\mathbb{R}) \rightarrow M_2(\mathbb{R})$ by setting $\lambda_A(B) = AB$, $\rho_A(B) = BA$, and $\kappa_A(B) = AB - BA$ for all $B \in M_2(\mathbb{R})$. (a) Confirm that λ_A , ρ_A , and κ_A are \mathbb{R} -linear. (b) Compute $[\lambda_A]_{X_1}$, $[\rho_A]_{X_1}$, and $[\kappa_A]_{X_1}$. (c) Compute $[\lambda_A]_{X_2}$, $[\rho_A]_{X_2}$, and $[\kappa_A]_{X_2}$.
17. Given $A \in M_{m,n}(F)$, define $\lambda_A : M_{n,p}(F) \rightarrow M_{m,p}(F)$ by $\lambda_A(B) = AB$ for $B \in M_{n,p}(F)$. Describe the entries of ${}_Y[\lambda_A]_X$, where X and Y are standard ordered bases for $M_{n,p}(F)$ and $M_{m,p}(F)$ (e.g., $X = (e_{11}, e_{12}, \dots, e_{1p}, e_{21}, \dots, e_{np})$).
18. Given $A \in M_{n,p}(F)$, define $\rho_A : M_{m,n}(F) \rightarrow M_{m,p}(F)$ by $\rho_A(B) = BA$ for $B \in M_{m,n}(F)$. Describe the entries of ${}_Y[\rho_A]_X$, where X and Y are standard ordered bases for $M_{m,n}(F)$ and $M_{m,p}(F)$.
19. Given $A \in M_n(F)$, define $\kappa_A : M_n(F) \rightarrow M_n(F)$ by $\kappa_A(B) = AB - BA$ for $B \in M_n(F)$. Let $X = (e_{11}, \dots, e_{1n}, e_{21}, \dots, e_{nn})$ be the standard ordered basis of $M_n(F)$. For each $i, j \in [n]$, compute $[\kappa_{e_{ij}}]_X$. Illustrate by writing out $[\kappa_{e_{23}}]_X$ for $n = 3$.
20. Show that each of the three bijections in (6.1) is F -linear. (This gives another proof that the map $T \mapsto {}_Y[T]_X$ is F -linear.)
21. Using only the definitions of spanning and linear independence (not matrices), prove that the maps T_{ij} defined in (6.4) form a basis of the F -vector space $L(V, W)$.
22. (a) Check that the map T defined in (6.2) is F -linear, well-defined, and sends x_j to w_j for $1 \leq j \leq n$. (b) Reprove (a) by showing that $T = L_{\mathbf{w}} \circ L_X^{-1}$, where $\mathbf{w} = (w_1, \dots, w_n) \in W^n$.
23. Define a map $M'_{X,Y} : L(V, W) \rightarrow M_{n,m}(F)$ by letting $M'_{X,Y}(T)$ be the $n \times m$ matrix whose j 'th row (for $1 \leq j \leq n$) is $[T(x_j)]_Y$ viewed as a row vector. (a) Prove or disprove: $M'_{X,Y}$ is a vector space isomorphism. (b) Let $V = W$, $X = Y$, and $m = n$. Prove or disprove: $M'_{X,X}$ is an F -algebra isomorphism.
24. Let \mathcal{B} be the set of pairs (X, Y) where X is an ordered basis for V and Y is an ordered basis for W , and let \mathcal{I} be the set of vector space isomorphisms $S : L(V, W) \rightarrow M_{m,n}(F)$. Define a map $\phi : \mathcal{B} \rightarrow \mathcal{I}$ by $\phi(X, Y) = M_{X,Y}$, where $M_{X,Y}(T) = {}_Y[T]_X$. (a) Prove ϕ is not injective in general. (b) Prove ϕ is not surjective in general.
25. Show that each list is an ordered basis of $P_{\leq 3}^*$. (a) (f_0, f_1, f_2, f_3) , where $f_i(c_0 + c_1t + c_2t^2 + c_3t^3) = c_i$ for $0 \leq i \leq 3$ and $c_j \in \mathbb{R}$. (b) (g_0, g_1, g_2, g_3) where $g_i(p) = (d/dt)^i(p)|_{t=0}$. (c) (h_0, h_1, h_2, h_3) where $h_i(p) = p(i)$.
26. Let $D : P_{\leq 3} \rightarrow P_{\leq 3}$ be the differentiation operator. For each basis X of $P_{\leq 3}^*$ found in Exercise 25, compute $[D^*]_X$.
27. Suppose $S \in L(V, W)$ and $U \in L(W, Z)$. Prove $(U \circ S)^* = S^* \circ U^*$. Translate this fact into an identity involving the matrix transpose.
28. Define $T : \mathbb{C} \rightarrow M_2(\mathbb{R})$ by $T(a + bi) = \begin{bmatrix} a & -b \\ b & a \end{bmatrix}$ for all $a, b \in \mathbb{R}$. Show that T is a one-to-one \mathbb{R} -algebra homomorphism.
29. Prove (6.6) by expressing $(S \circ T)(x_j)$ as a specific linear combination of vectors in Z and showing that the coefficient of z_i in this combination is $(BA)(i, j)$.
30. Suppose $T : V \rightarrow W$ is an F -linear map. You are given ${}_Y[T]_X$ but do not know any other specific information about T , X , or Y . (a) Explain how to use

- matrix algorithms to compute a basis for $\ker(T)$. (b) Explain how to use matrix algorithms to compute a basis for $\text{img}(T)$.
31. Let $X_4 = ((1, 1, 1), (1, 0, -1), (1, 0, 1))$, which is an ordered basis of \mathbb{R}^3 . Compute ${}_{X_4}[\text{id}]_{X_4}$ and ${}_{X_i}[\text{id}]_{X_i}$ for $i = 1, 2, 3$.
 32. For $V = P_{\leq 3}$, compute the transition matrix: (a) from X_3 to X_1 ; (b) from X_1 to X_3 ; (c) from X_3 to X_2 ; (d) from X_2 to X_3 .
 33. Compute transition matrices between all pairs of ordered bases for \mathbb{C} listed in Table 6.1.
 34. Compute transition matrices between all pairs of ordered bases for $M_2(\mathbb{R})$ listed in Table 6.1.
 35. (a) Given $[v]_{X_3} = (5, -1, -1)$ in \mathbb{R}^3 , find $[v]_{X_2}$. (b) Given $[v]_{X_2} = (2, 1, 1, -2)$ in $P_{\leq 3}$, find $[v]_{X_3}$. (c) Given $[v]_{X_2} = (4, -5, 0, 1)$ in $M_2(\mathbb{R})$, find $[v]_{X_1}$ and $[v]_{X_3}$.
 36. Let $P_{< n} = \{f \in \mathbb{R}[t] : f = 0 \text{ or } \deg(f) < n\}$. $P_{< n}$ has ordered bases $X = (1, t, t^2, \dots, t^{n-1})$ and $Y = (1, (t+c), (t+c)^2, \dots, (t+c)^{n-1})$ for fixed $c \in \mathbb{R}$. (a) Find ${}_{X}[\text{id}]_{Y}$. (b) Find ${}_{Y}[\text{id}]_{X}$.
 37. Let V consist of polynomials in $\mathbb{R}[t]$ of degree at most 4. Let $X = (1, t, t^2, t^3, t^4)$ and $Y = (1, t, t(t-1), t(t-1)(t-2), t(t-1)(t-2)(t-3))$. (a) Find ${}_{X}[\text{id}]_{Y}$. (b) Find ${}_{Y}[\text{id}]_{X}$.
 38. Prove that for a fixed ordered basis Y of V , every invertible matrix in $M_n(F)$ has the form ${}_{X}[\text{id}]_{Y}$ for a unique ordered basis X of V .
 39. Let $A = \begin{bmatrix} 2 & 1 & -1 \\ 2 & 3 & 0 \\ 4 & 1 & -2 \end{bmatrix}$. (a) Find the unique ordered basis Z of \mathbb{R}^3 such that ${}_{X_3}[\text{id}]_Z = A$. (b) Find the unique ordered basis Y of \mathbb{R}^3 such that ${}_{Y}[\text{id}]_{X_3} = A$.
 40. Verify that each relation is an equivalence relation. (a) column-equivalence on $M_{m,n}(F)$; (b) row-equivalence on $M_{m,n}(F)$; (c) row/column-equivalence on $M_{m,n}(F)$; (d) similarity on $M_n(F)$; (e) congruence on $M_n(F)$; (f) orthogonal similarity on $M_n(\mathbb{R})$; (g) unitary similarity on $M_n(\mathbb{C})$.
 41. (a) Find all the similarity equivalence classes in $M_2(\mathbb{Z}_2)$. (b) Find all the congruence equivalence classes in $M_2(\mathbb{Z}_2)$.
 42. Find all similarity equivalence classes in $M_n(F)$ of size one.
 43. Let $D : P_{\leq 3} \rightarrow P_{\leq 3}$ be given by $D(f) = f'$ for $f \in P_{\leq 3}$. (a) Find ordered bases X and Y of $P_{\leq 3}$ such that ${}_{Y}[D]_X$ is diagonal with zeroes and ones on the diagonal. (b) Find an ordered basis X of $P_{\leq 3}$ such that $[D]_X$ is diagonal, or explain why this cannot be done.
 44. Let $A \in M_n(F)$ be an invertible matrix. (a) Explicitly describe ordered bases X and Y of F^n such that ${}_{Y}[L^A]_X = I_n$. (b) Prove: if $[L^A]_X = I_n$ for some ordered basis X , then $A = I_n$.
 45. Prove the projection theorem stated in the last paragraph of §6.11 without using matrices.
 46. Let $A \in M_n(F)$ and $P \in \text{GL}_n(F)$. Prove $x \in F^n$ is an eigenvector of A with associated eigenvalue c iff $P^{-1}x$ is an eigenvector of $P^{-1}AP$ with associated eigenvalue c . Conclude that similar matrices have the same eigenvalues.
 47. (a) Show that for all $A \in M_n(F)$ and $c \in F$, A is diagonalizable iff $A + cI_n$ is diagonalizable. (b) Show that a strictly upper-triangular nonzero matrix A is not

- diagonalizable. (c) Show that if $c \in F$ is the only eigenvalue of A , then $A = cI_n$ or A is not diagonalizable.
48. (a) Show that if $A \in M_n(F)$ has no eigenvalues in F , then A is not triangulable.
 (b) Show that $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ is not triangulable in $M_2(\mathbb{R})$, but A is diagonalizable in $M_2(\mathbb{C})$. (c) Give an example of a field F and $A \in M_n(F)$ such that A has an eigenvalue in F , but A is not triangulable.
49. Decide if each matrix A is diagonalizable. If so, find a diagonal D and an invertible P with $P^{-1}AP = D$. (a) $A = \frac{1}{7} \begin{bmatrix} 32 & 12 & -6 \\ 9 & 20 & -3 \\ 9 & 6 & 11 \end{bmatrix} \in M_3(\mathbb{R})$.
 (b) $A = \begin{bmatrix} 1 & 1 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & -1 \end{bmatrix} \in M_3(\mathbb{R})$. (c) the matrix in (b), viewed as an element of $M_3(\mathbb{C})$. (d) $A = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix} \in M_4(\mathbb{Q})$. (e) The matrix in (d), viewed as an element of $M_4(\mathbb{Z}_2)$.
50. Given a complete flag of subspaces (V_0, V_1, \dots, V_n) of V , let x_i be any element of $V_i \sim V_{i-1}$, for $1 \leq i \leq n$. Prove $X = (x_1, \dots, x_n)$ is an ordered basis of V .
51. Let (x_1, \dots, x_n) be an ordered basis for V such that (x_1, \dots, x_k) is an ordered basis for a subspace W . Prove $(x_{k+1} + W, \dots, x_n + W)$ is an ordered basis for V/W .
52. (a) Give an example of a map $T \in L(\mathbb{R}^2)$ such that the only T -invariant subspaces of \mathbb{R}^2 are $\{0\}$ and \mathbb{R}^2 . (b) Give an example of a map $T \in L(\mathbb{R}^2)$ such that \mathbb{R}^2 has exactly four T -invariant subspaces. (c) Does there exist $T \in L(\mathbb{R}^2)$ such that T has exactly three T -invariant subspaces? (d) Give an example of $T \in L(\mathbb{R}^3)$ such that \mathbb{R}^3 has infinitely many T -invariant subspaces, but T is not a scalar multiple of id.
53. Given $i_1, i_2, \dots, i_b \in \mathbb{N}^+$ with sum n , let
- $$I_1 = \{1, 2, \dots, i_1\}, I_2 = \{i_1 + 1, i_1 + 2, \dots, i_1 + i_2\}, \dots, \\ I_b = \{i_1 + \dots + i_{b-1} + 1, \dots, n\}.$$
- Call a matrix $A \in M_n(F)$ *block-triangular* with block sizes i_1, i_2, \dots, i_b iff for all $i \in I_r$ and $j \in I_s$ with $r > s$, $A(i, j) = 0$. Define a *partial flag* of type (i_1, \dots, i_b) to be a chain of subspaces $V_1 \subseteq V_2 \subseteq \dots \subseteq V_b$ with $\dim(V_j) = i_1 + i_2 + \dots + i_j$ for $1 \leq j \leq b$. Show that there exists an ordered basis X of V such that $[T]_X$ is block-triangular with block sizes i_1, i_2, \dots, i_b iff V has a partial flag of type (i_1, \dots, i_b) stabilized by T .
54. Define I_1, \dots, I_b as in Exercise 53. Call a matrix $A \in M_n(F)$ *block-diagonal* with block sizes i_1, \dots, i_b iff for all $i \in I_r$ and $j \in I_s$ with $r \neq s$, $A(i, j) = 0$. Find and prove an abstract criterion on $T \in L(V)$ that is equivalent to the existence of an ordered basis X of V such that $[T]_X$ is block-diagonal with block sizes i_1, \dots, i_b .
55. Give an example of a 4×4 real matrix A with all entries nonzero such that A is similar to a block-triangular matrix, but A is not triangulable.

56. Let $A = \begin{bmatrix} 1/2 & -1/2 & -1/2 & 1/2 \\ 1/2 & 5/2 & 5/2 & 1/2 \\ -1/2 & -3/2 & -3/2 & -1/2 \\ 1/2 & 1/2 & 1/2 & 1/2 \end{bmatrix}$. (a) Verify that A is idempotent.
 (b) Find an ordered basis X of \mathbb{R}^4 with $[L^A]_X = I_{k,4}$ for some k . (c) Find an invertible $P \in M_4(\mathbb{R})$ with $P^{-1}AP = I_{k,4}$.
57. Given an idempotent matrix $A \in M_n(F)$, describe an algorithm for finding k and $P \in \text{GL}_n(F)$ with $P^{-1}AP = I_{k,n}$.
58. *Projections of a Direct Sum.* We say that V is the *direct sum* of subspaces W_1, \dots, W_k , denoted $V = W_1 \oplus \dots \oplus W_k$, iff every $v \in V$ can be written uniquely in the form $v = w_1 + \dots + w_k$ with $w_i \in W_i$. (a) Given $V = W_1 \oplus \dots \oplus W_k$, define maps $P_1, \dots, P_k : V \rightarrow V$ by setting $P_i(v) = w_i$, where $v = w_1 + \dots + w_k$ as above. Show that each P_i is linear and idempotent with image W_i , $P_iP_j = 0$ for all $i \neq j$, and $P_1 + \dots + P_k = \text{id}_V$. What is $\ker(P_i)$? (b) Let X_1, \dots, X_k be ordered bases for W_1, \dots, W_k . Let X be the concatenation of X_1, \dots, X_k . Show that X is an ordered basis of V , and find $[P_i]_X$ for each i .
59. Suppose Q_1, \dots, Q_k are idempotent elements of $L(V)$ such that $Q_iQ_j = 0$ for all $i \neq j$ and $Q_1 + \dots + Q_k = \text{id}_V$. (a) Writing $W_i = \text{img}(Q_i)$ for $1 \leq i \leq k$, show that $V = W_1 \oplus \dots \oplus W_k$. (b) Defining P_i as in Exercise 58, show that $P_i = Q_i$ for all i .
60. Suppose $V = W \oplus Z$, and let $P = P_{W,Z}$ be the associated projection. (a) Prove: W is T -invariant iff $PTP = TP$ in $L(V)$. (b) Prove: W and Z are both T -invariant iff $PT = TP$ in $L(V)$.
61. (a) Check that $\text{BL}(V)$ is a subspace of the vector space of all functions from $V \times V$ to F . (b) Check that the map N_X defined in §6.18 is F -linear. (c) Prove (6.8) by induction. (d) Complete the proof that the map B defined at the end of §6.18 is F -bilinear and satisfies $B(x_i, x_j) = A(i, j)$.
62. Check that each map B is F -bilinear, and then compute $[B]_X$ for the indicated ordered basis X . (a) $B\left(\begin{bmatrix} a \\ c \end{bmatrix}, \begin{bmatrix} b \\ d \end{bmatrix}\right) = \det \begin{bmatrix} a & b \\ c & d \end{bmatrix}$, using $X = (e_1, e_2)$.
 (b) For $f, g \in P_{\leq 3}$, $B(f, g) = \int_0^1 f(t)g(t) dt$, using $X = (1, t, t^2, t^3)$. (c) For $z, w \in \mathbb{C}$, $B(z, w) =$ the real part of zw , using $X = (1, i)$.
63. Use your answers to each part of Exercise 62 to find $[B]_Y$ for the new ordered basis Y . (a) $Y = ((1, 2), (3, 4))$. (b) $Y = (t + 3, 2t^2 - 4, t^3 - t^2, t^3 + t^2)$. (c) $Y = (3 + 4i, 2 - i)$.
64. For $V = \mathbb{R}^3$, we are given $B \in \text{BL}(V)$ with $[B]_{X_2} = \begin{bmatrix} 2 & 1 & 0 \\ -1 & 3 & 1 \\ 0 & 1 & -1 \end{bmatrix}$.
 (a) Find $B((1, 0, -4), (1, 3, 7))$. (b) Find $[B]_{X_1}$. (c) Find $[B]_{X_3}$.
65. Given $B \in \text{BL}(V)$ and $v, w \in V$, show that $[B(v, w)] = ([v]_X)^T [B]_X [w]_X$.
66. (a) How does the matrix $[B]_X$ change if we switch the positions of x_i and x_j in X ? (b) How does the matrix $[B]_X$ change if we multiply x_i by a nonzero $c \in F$?
67. Given $B \in \text{BL}(V)$ and $v, w \in V$, define $L_v(w) = B(v, w)$ and $R_v(w) = B(w, v)$.
 (a) Prove: for all $v \in V$, L_v and R_v lie in V^* . (b) How is the row vector ${}_{(1_F)}[L_{x_i}]_X$ related to the matrix $[B]_X$? How is the row vector ${}_{(1_F)}[R_{x_j}]_X$ related to the matrix $[B]_X$? (c) More generally, describe how to use $[B]_X$ to compute ${}_{(1_F)}[L_v]_X$ and ${}_{(1_F)}[R_v]_X$ for any $v \in V$.

68. Let \mathcal{B} be the set of ordered bases for V , and let \mathcal{I} be the set of vector space isomorphisms $S : \text{BL}(V) \rightarrow M_n(F)$. Define $\phi : \mathcal{B} \rightarrow \mathcal{I}$ by $\phi(X) = N_X$, where $N_X(B) = [B]_X$ for $B \in \text{BL}(V)$. Prove or disprove: ϕ is a bijection.
69. (a) For a real inner product space V , check that $O(V)$ is a subgroup of $\text{GL}(V)$.
 (b) For a complex inner product space V , check that $U(V)$ is a subgroup of $\text{GL}(V)$. (c) Check that $O_n(\mathbb{R})$ is a subgroup of $\text{GL}_n(\mathbb{R})$. (d) Check that $U_n(\mathbb{C})$ is a subgroup of $\text{GL}_n(\mathbb{C})$.
70. Fix an orthonormal basis X of a real inner product space V . (a) Show that every orthogonal matrix $P \in O_n(\mathbb{R})$ has the form ${}_Y[\text{id}]_X$ for a unique orthonormal basis Y of V . (b) Show that every orthogonal matrix $P \in O_n(\mathbb{R})$ has the form ${}_X[\text{id}]_Y$ for a unique orthonormal basis Y of V . (c) Prove analogues of (a) and (b) for complex inner product spaces.
71. Let $v = (v_1, v_2, v_3)$ and $w = (w_1, w_2, w_3)$ be vectors in \mathbb{R}^3 with $\|v\| = 1 = \|w\|$ and $v \bullet w = 0$. Find two vectors $z \in \mathbb{R}^3$ such that (v, w, z) is an orthonormal basis of \mathbb{R}^3 .
72. (a) Show that $P_{\leq 3}$ becomes a real inner product space if we define $\langle f, g \rangle = \int_0^1 f(t)g(t) dt$ for $f, g \in P_{\leq 3}$. (b) Find an orthonormal basis (f_0, f_1, f_2, f_3) of $P_{\leq 3}$ such that $\deg(f_i) = i$ for $0 \leq i \leq 3$.
73. (a) For each $t \in \mathbb{R}$, show that $A_t = \begin{bmatrix} \cos t & -\sin t \\ \sin t & \cos t \end{bmatrix}$ is an orthogonal matrix.
 (b) Show that every orthogonal matrix in $O_2(\mathbb{R})$ has the form A_t or $A_t \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$ for some $t \in \mathbb{R}$.
74. Let \mathcal{B} be the set of orthonormal bases for a real inner product space V , and let \mathcal{I} be the set of orthogonal vector space isomorphisms $S : F^n \rightarrow V$. Define $\phi : \mathcal{B} \rightarrow \mathcal{I}$ by $\phi(X) = L_X$ for all $X \in \mathcal{B}$. Prove or disprove: ϕ is a bijection.
75. Let \mathcal{B} be the set of ordered bases for V , and let \mathcal{I} be the set of F -algebra isomorphisms $S : L(V) \rightarrow M_n(F)$. Define $\phi : \mathcal{B} \rightarrow \mathcal{I}$ by $\phi(X) = M_X$, where $M_X(T) = [T]_X$ for $T \in L(V)$. Determine (with proof) whether ϕ is injective or surjective.

Part III

Matrices with Special Structure

This page intentionally left blank

Hermitian, Positive Definite, Unitary, and Normal Matrices

We begin by recalling some fundamental facts about complex numbers.

- For every complex number z , there exist unique real numbers x and y such that $z = x + iy$. We call x the *real part* of z and y the *imaginary part* of z , writing $x = \operatorname{Re}(z)$ and $y = \operatorname{Im}(z)$. The formula $z = x + iy$ is called the *Cartesian decomposition* of the complex number z .
- For each complex number $z = x + iy$ with $x, y \in \mathbb{R}$, the *complex conjugate* of z is $\bar{z} = x - iy$. For all $z, w \in \mathbb{C}$, the following properties hold: $\bar{\bar{z}} = z$; $\bar{z + w} = \bar{z} + \bar{w}$; $\bar{zw} = \bar{z} \cdot \bar{w} = \bar{w} \cdot \bar{z}$; $\operatorname{Re}(z) = (z + \bar{z})/2$; $\operatorname{Im}(z) = (z - \bar{z})/2i$; z is real iff $\bar{z} = z$; z is pure imaginary (meaning $\operatorname{Re}(z) = 0$) iff $\bar{z} = -z$.
- The *magnitude* of a complex number $z = x + iy$ is the nonnegative real number $|z| = \sqrt{x^2 + y^2} = (\bar{z}z)^{1/2}$. Every nonzero complex number can be written uniquely in the form $z = rw$, where r is a positive real number and w is a complex number of magnitude 1. We often write $w = e^{i\theta} = \cos \theta + i \sin \theta$ for some (non-unique) real number θ ; then $z = re^{i\theta}$ is called the *polar decomposition* of z .

Now consider the set $M_n(\mathbb{C})$ of all $n \times n$ matrices with complex entries. If $n = 1$, we can identify $M_n(\mathbb{C})$ with \mathbb{C} . Our goal in this chapter is to build an analogy between \mathbb{C} and $M_n(\mathbb{C})$ for $n > 1$. We will define matrix versions of the concepts of real numbers, positive numbers, pure imaginary numbers, complex numbers of magnitude 1, complex conjugation, the Cartesian decomposition of a complex number, and the polar decomposition of a complex number (among other things). This approach allows us to unify a diverse collection of fundamental results in matrix theory. The analogy with \mathbb{C} also provides some motivation for the introduction of certain special classes of matrices (namely unitary, Hermitian, positive definite, and normal matrices) that have many remarkable properties.

7.1 Conjugate-Transpose of a Matrix

The first step towards implementing the analogy between complex numbers and complex-valued matrices is to introduce a matrix version of complex conjugation. One possibility is to consider the operation on matrices that replaces each entry of a matrix by its complex conjugate. However, it turns out to be much more fruitful to combine this operation with the transpose map that interchanges the rows and columns of a matrix.

Formally, suppose A is an $m \times n$ matrix with complex entries $A(i, j)$, where $1 \leq i \leq m$ and $1 \leq j \leq n$. Define the *conjugate-transpose* of A to be the $n \times m$ matrix A^* such that

$$A^*(i, j) = \overline{A(j, i)} \quad (1 \leq i \leq n, 1 \leq j \leq m).$$

For example, given $A = \begin{bmatrix} 2+i & 3 & 5i \\ 1-2i & 0 & -1+3i \end{bmatrix}$, we compute $A^* = \begin{bmatrix} 2-i & 1+2i \\ 3 & 0 \\ -5i & -1-3i \end{bmatrix}$.

In the case of a 1×1 matrix, the conjugate-transpose operation reduces to the complex conjugation operation on \mathbb{C} .

The conjugate-transpose operation shares many of the algebraic properties of complex conjugation. Suppose $A, B \in M_{m,n}(\mathbb{C})$ and $z \in \mathbb{C}$ is a complex scalar. First, $(A+B)^* = A^* + B^*$, since for all $i \in [n]$ and $j \in [m]$,

$$(A+B)^*(i,j) = \overline{(A+B)(j,i)} = \overline{A(j,i) + B(j,i)} = \overline{A(j,i)} + \overline{B(j,i)} = A^*(i,j) + B^*(i,j).$$

Second, $(zA)^* = \overline{z}(A^*)$, since

$$(zA)^*(i,j) = \overline{zA(j,i)} = \overline{z}\overline{A(j,i)} = \overline{z}(A^*(i,j)) = (\overline{z}A^*)(i,j).$$

In particular, $(xA)^* = x(A^*)$ if x is a *real* scalar. Third, if $D \in M_{n,p}(\mathbb{C})$, then $(AD)^* = D^*A^*$. To prove this, first note that both sides are $p \times m$ matrices. For $1 \leq i \leq p$ and $1 \leq j \leq m$, compute

$$\begin{aligned} (AD)^*(i,j) &= \overline{AD(j,i)} = \overline{\sum_{k=1}^n A(j,k)D(k,i)} \\ &= \sum_{k=1}^n \overline{A(j,k)} \cdot \overline{D(k,i)} \\ &= \sum_{k=1}^n D^*(i,k)A^*(k,j) = (D^*A^*)(i,j). \end{aligned}$$

Fourth, $(A^*)^* = A$, since for $i \in [m]$ and $j \in [n]$,

$$(A^*)^*(i,j) = \overline{A^*(j,i)} = \overline{\overline{A(i,j)}} = A(i,j).$$

Fifth, we see by induction on $s \in \mathbb{N}^+$ that

$$(z_1A_1 + \cdots + z_sA_s)^* = \overline{z_1}A_1^* + \cdots + \overline{z_s}A_s^*; \quad (B_1B_2 \cdots B_s)^* = B_s^* \cdots B_2^*B_1^*$$

whenever these expressions are defined. Sixth, if $A \in M_n(\mathbb{C})$ is invertible, then A^* is invertible and $(A^{-1})^* = (A^*)^{-1}$. To see this, write $B = A^{-1}$, and apply the conjugate-transpose to the identities $AB = I = BA$ to obtain $B^*A^* = I^* = I = A^*B^*$. Thus, $B^* = (A^{-1})^*$ is the two-sided matrix inverse for A^* , as needed. Seventh, for $A \in M_n(\mathbb{C})$, $\text{tr}(A^*) = \overline{\text{tr}(A)}$, where tr denotes the trace (the sum of the diagonal entries of A). This follows since

$$\text{tr}(A^*) = \sum_{i=1}^n A^*(i,i) = \sum_{i=1}^n \overline{A(i,i)} = \overline{\sum_{i=1}^n A(i,i)} = \overline{\text{tr}(A)}.$$

Eighth, for $A \in M_n(\mathbb{C})$, $\det(A^*) = \overline{\det(A)}$, since

$$\begin{aligned} \det(A^*) &= \sum_{f \in S_n} \text{sgn}(f)A^*(f(1),1)A^*(f(2),2) \cdots A^*(f(n),n) \\ &= \sum_{f \in S_n} \text{sgn}(f)\overline{A(1,f(1))} \cdot \overline{A(2,f(2))} \cdots \overline{A(n,f(n))} \\ &= \overline{\sum_{f \in S_n} \text{sgn}(f)A^T(f(1),1)A^T(f(2),2) \cdots A^T(f(n),n)} \\ &= \overline{\det(A^T)} = \overline{\det(A)}. \end{aligned}$$

(See Chapter 5 for further discussion of determinants.)

The conjugate-transpose operation is related to the standard inner product on \mathbb{C}^n . If $v = (v_1, \dots, v_n)$ and $w = (w_1, \dots, w_n)$ belong to \mathbb{C}^n , their *inner product* is defined to be $\langle v, w \rangle = \sum_{i=1}^n \overline{w_i} v_i$. If we identify \mathbb{C}^n with the space of column vectors $M_{n,1}(\mathbb{C})$, and if we identify 1×1 matrices with complex numbers, then the definition of matrix multiplication shows that $\langle v, w \rangle = w^* v$. Observe that $v^* v = \langle v, v \rangle = \sum_{i=1}^n \overline{v_i} v_i = \sum_{i=1}^n |v_i|^2$ is a nonnegative real number, which is zero iff $v = 0$. We define the *Euclidean norm* of v to be $\|v\| = (v^* v)^{1/2}$ for all $v \in \mathbb{C}^n$. If we restrict v to lie in \mathbb{R}^n , the complex conjugates disappear and we recover the standard inner product (the dot product) and the Euclidean norm on \mathbb{R}^n . Note the following *adjoint property* of the matrix A^* : for any $A \in M_n(\mathbb{C})$ and any $v, w \in \mathbb{C}^n$,

$$\langle Av, w \rangle = \langle v, A^* w \rangle.$$

This follows since both sides are equal to $w^* Av$.

Given a matrix $A \in M_n(\mathbb{C})$, we can use the conjugate-transpose operation to define the *quadratic form* on \mathbb{C}^n associated with A . This quadratic form is the function $Q : \mathbb{C}^n \rightarrow \mathbb{C}$ defined by $Q(v) = v^* Av$ for $v \in \mathbb{C}^n$. To get a formula for $Q(v)$ in terms of the components of v , write $v = (v_1, \dots, v_n)$ with each $v_k \in \mathbb{C}$. The k 'th component of the column vector Av is $\sum_{j=1}^n A(k, j)v_j$. Multiplying this column vector on the left by the row vector v^* , we see that

$$Q(v) = \sum_{k=1}^n \sum_{j=1}^n \overline{v_k} A(k, j)v_j \quad (v \in \mathbb{C}^n). \quad (7.1)$$

For example, if $A = \begin{bmatrix} 1 & 2-i \\ 3i & 0 \end{bmatrix}$ and $v = (3+i, -4i)$, then

$$Q(v) = (3-i)(3+i) + (3-i)(2-i)(-4i) + (4i)(3i)(3+i) + (4i)0(-4i) = -46 - 32i.$$

7.2 Hermitian Matrices

Recall that a complex number z is real iff $\bar{z} = z$. This suggests the following definition: a complex-valued matrix A is *Hermitian* iff $A^* = A$. This equality is only possible if A is a square matrix. We see from the definition of A^* that $A \in M_n(\mathbb{C})$ is Hermitian iff $A(i, j) = \overline{A(j, i)}$ for all $1 \leq i, j \leq n$. For example,

$$A = \begin{bmatrix} 2 & 1+8i & 0 & i \\ 1-8i & -1 & 3-5i & 7+i/2 \\ 0 & 3+5i & 0 & 1 \\ -i & 7-i/2 & 1 & \pi \end{bmatrix}$$

is a Hermitian matrix. Observe that a real-valued matrix A is Hermitian iff A is symmetric (i.e., $A^T = A$). Furthermore, for any Hermitian matrix A , the main diagonal entries of A must all be real. This follows from the identity $A(i, i) = \overline{A(i, i)}$. A 1×1 complex matrix is Hermitian iff the sole entry of that matrix is real. The reader should think of Hermitian matrices as being the matrix analogues of real numbers. This analogy will be strengthened by results proved later in this section.

Suppose A and B are $n \times n$ Hermitian matrices and x is a *real* scalar. Then $A + B$ and xA are Hermitian, since $(A + B)^* = A^* + B^* = A + B$ and $(xA)^* = x(A^*) = xA$. More generally, any real linear combination of Hermitian matrices is also Hermitian. On

the other hand, AB need not be Hermitian in general, since $(AB)^* = B^*A^* = BA$. This equation shows that $(AB)^* = AB$ holds iff $AB = BA$. Thus, the product of two *commuting* Hermitian matrices is also Hermitian. Finally, if A is Hermitian and B is arbitrary, then B^*AB is also Hermitian. For, $(B^*AB)^* = B^*A^*(B^*)^* = B^*AB$.

If A is Hermitian and $v, w \in \mathbb{C}^n$, we know that $\langle Av, w \rangle = \langle v, A^*w \rangle = \langle v, Aw \rangle$. This fact is sometimes expressed by saying that a Hermitian matrix is *self-adjoint*. Next, consider the quadratic form Q associated to a Hermitian matrix A . We claim that, for all $v \in \mathbb{C}^n$, $Q(v) = v^*Av$ is *real*. To confirm this, we need only check that $\overline{Q(v)} = Q(v)$. Now, since v^*Av is a 1×1 matrix and A is Hermitian,

$$\overline{Q(v)} = \overline{v^*Av} = (v^*Av)^* = v^*A^*(v^*)^* = v^*Av = Q(v).$$

Another key fact about Hermitian matrices is that *all eigenvalues of a Hermitian matrix A must be real*. To prove this, let $\lambda \in \mathbb{C}$ be an eigenvalue of A with an associated nonzero eigenvector $v = (v_1, \dots, v_n) \in \mathbb{C}^n$. Direct calculation shows that $v^*v = \sum_{i=1}^n |v_i|^2$ is a positive real number. We have just seen that v^*Av is also real. But $v^*Av = v^*(\lambda v) = \lambda(v^*v)$. Thus, $\lambda = (v^*Av)/(v^*v)$ is a quotient of two real numbers, and is therefore real. We remark that a real-valued matrix $B \in M_n(\mathbb{R})$ need not have all real eigenvalues. For example, the eigenvalues of $B = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$ are readily found to be i and $-i$. This gives more evidence that the Hermitian property (as opposed to the more naive requirement of having all real entries) is the “right” generalization of “real-number-hood” to matrices. Note here that iB is Hermitian although every entry of iB is pure imaginary!

The property of having a real-valued quadratic form actually characterizes the class of Hermitian matrices. More precisely, if $A \in M_n(\mathbb{C})$ is a matrix such that $Q(v) = v^*Av \in \mathbb{R}$ for all $v \in \mathbb{C}^n$, then A must be Hermitian. To prove this, consider the unit vectors $e_k \in \mathbb{C}^n$, which have a 1 in position k and zeroes elsewhere. Using (7.1), we see that $Q(e_k) = A(k, k)$ must be real for $1 \leq k \leq n$. Next, for all $j \neq k$, we see similarly that $Q(e_j + e_k) = A(j, j) + A(j, k) + A(k, j) + A(k, k)$ is real, and therefore $A(j, k) + A(k, j) \in \mathbb{R}$ for all $j \neq k$. Finally, since we know $Q(e_j + ie_k) = A(j, j) + iA(j, k) + \bar{i}A(k, j) + A(k, k)$ is real, it follows that $i(A(j, k) - A(k, j))$ is real for all $j \neq k$. Fix $j \neq k$, and write $A(j, k) = a + ib$ and $A(k, j) = c + id$ with $a, b, c, d \in \mathbb{R}$. The preceding remarks say that $b + d = 0$ and $a - c = 0$, so $A(k, j) = c + id = a - ib = \overline{A(j, k)}$ for all $j \neq k$. We also have $A(j, j) = \overline{A(j, j)}$ for all j , since $A(j, j)$ is real. Thus, $A = A^*$ as claimed.

On the other hand, a matrix $A \in M_n(\mathbb{R})$ with all real eigenvalues need not be Hermitian. For instance, $A = \begin{bmatrix} 1 & 2 \\ 0 & 3 \end{bmatrix}$ has eigenvalues 1 and 3 (the entries on the main diagonal), but $A \neq A^*$. Soon, we will introduce a new matrix property (called normality) and prove that a *normal* matrix with all real eigenvalues must be Hermitian.

7.3 Hermitian Decomposition of a Matrix

Now consider an arbitrary matrix $A \in M_n(\mathbb{C})$. The analogy with \mathbb{C} suggests the following theorem: *there exist unique Hermitian matrices $X, Y \in M_n(\mathbb{C})$ such that $A = X + iY$* . To prove existence, recall how complex conjugation can be used to recover the real and imaginary parts of a complex number $z = x + iy$: $x = (z + \bar{z})/2$ and $y = (z - \bar{z})/2i$. Given the matrix A , we are therefore led to consider the matrices $X = (A + A^*)/2$ and $Y = (A - A^*)/2i$. We have $X^* = (A^* + A)/2 = X$ and $Y^* = \overline{\frac{1}{2i}}(A^* - A) = \frac{1}{-2i}(A^* - A) = Y$, so X and Y

are Hermitian. A quick calculation confirms that $X + iY = A$. Next, to prove uniqueness, suppose that $A = X + iY = U + iV$ where X, Y, U, V are all Hermitian. We must prove $X = U$ and $Y = V$. Observe that $X - U = i(V - Y)$. Applying the conjugate-transpose operation to both sides of this matrix identity, we get $X - U = -i(V - Y)$. Thus, $i(V - Y) = -i(V - Y)$, which implies $V - Y = 0$ and hence $Y = V$. Then $X - U = i(V - Y) = 0$ implies $X = U$, as needed.

The preceding theorem can also be rephrased in terms of skew-Hermitian matrices. A complex-valued matrix A is called *skew-Hermitian* iff $A^* = -A$. (This concept is the analogue of a pure imaginary number in \mathbb{C} .) One checks readily that a real linear combination of skew-Hermitian matrices is also skew-Hermitian. The zero matrix is the only matrix that is both Hermitian and skew-Hermitian. Moreover, a matrix B is Hermitian iff iB is skew-Hermitian, since $(iB)^* = -i(B^*)$. Combining this fact with the previous theorem, we see that any matrix $A \in M_n(\mathbb{C})$ can be written uniquely in the form $A = H + S$, where H is Hermitian and S is skew-Hermitian. More precisely, we have $H = X = (A + A^*)/2$ and $S = iY = (A - A^*)/2i$. H and S are respectively called the *Hermitian part of A* and the *skew-Hermitian part of A* . The decomposition $A = X + iY = H + S$ is called the *Hermitian decomposition* or *Cartesian decomposition* of the matrix A .

We pause to point out a seemingly innocuous property of the Cartesian decomposition of a complex number. Suppose we write $z \in \mathbb{C}$ in the form $z = x + iy$, where $x, y \in \mathbb{R}$. Then $xy = yx$, since multiplication of real numbers is commutative. The analogous property for complex matrices is much more subtle — indeed, it need not even be true in general! Suppose $A = X + iY$ is the Cartesian decomposition of $A \in M_n(\mathbb{C})$. Recall that $X = (A + A^*)/2$ and $Y = (A - A^*)/2i$. A calculation with the distributive law reveals that

$$XY = \frac{1}{4i}(AA + A^*A - AA^* - A^*A^*); \quad YX = \frac{1}{4i}(AA - A^*A + AA^* - A^*A^*).$$

Comparing these expressions, we see that $XY = YX$ holds iff $A^*A - AA^* = -A^*A + AA^*$ iff $2A^*A = 2AA^*$ iff $A^*A = AA^*$. The latter condition does not hold for every matrix; for example, if $A = \begin{bmatrix} 1 & 2 \\ 0 & 3 \end{bmatrix}$, then $A^* = \begin{bmatrix} 1 & 0 \\ 2 & 3 \end{bmatrix}$ and

$$AA^* = \begin{bmatrix} 5 & 6 \\ 6 & 9 \end{bmatrix}; \quad A^*A = \begin{bmatrix} 1 & 2 \\ 2 & 13 \end{bmatrix}.$$

We say that a matrix $A \in M_n(\mathbb{C})$ is a *normal matrix* iff $A^*A = AA^*$. We shall soon see that normal matrices have very nice properties compared to general complex matrices.

7.4 Positive Definite Matrices

Our next task is to find the matrix analogues of positive and negative real numbers. To this end, recall that a matrix $A \in M_n(\mathbb{C})$ is Hermitian iff $A = A^*$ iff $Q(v) = v^*Av$ is *real* for all $v \in \mathbb{C}^n$ (§7.2). This suggests the following definition: we say a matrix $A \in M_n(\mathbb{C})$ is *positive definite* iff v^*Av is a positive real number for all nonzero $v \in \mathbb{C}^n$. We say A is *positive semidefinite* iff v^*Av is a nonnegative real number for all nonzero $v \in \mathbb{C}^n$. Similarly, A is *negative definite* (resp. *negative semidefinite*) iff v^*Av is a negative real number (resp. a nonpositive real number) for all nonzero $v \in \mathbb{C}^n$. Observe that positive definite matrices must be Hermitian, as one would expect (since positive real numbers are real); similarly for positive semidefinite, negative definite, and negative semidefinite matrices.

Suppose $A, B \in M_n(\mathbb{C})$ are positive definite and $c > 0$ is a positive real scalar. Then $A + B$ is positive definite, since $v^*(A + B)v = v^*Av + v^*Bv$ is a sum of two positive real numbers for all nonzero $v \in \mathbb{C}^n$; and cA is positive definite, since $v^*(cA)v = c(v^*Av)$ is a product of two positive real numbers. More generally, any linear combination of positive definite matrices with positive real coefficients is also positive definite. Similarly, any linear combination of negative definite matrices with *positive* real coefficients is also negative definite. Moreover, A is positive definite iff $-A$ is negative definite. Analogous properties hold for semidefinite matrices.

Another helpful fact about positive definite matrices is that *all eigenvalues of a positive definite matrix A are strictly positive real numbers*. To prove this, let $\lambda \in \mathbb{C}$ be an eigenvalue of A , and let v be an associated nonzero eigenvector. Recall that $v^*v = \sum_{i=1}^n |v_i|^2$ is a positive real number. By positive definiteness, so too is $v^*Av = v^*(\lambda v) = \lambda(v^*v)$. But then $\lambda = (v^*Av)/(v^*v)$ is also a positive real number, being the ratio of two positive real numbers. Exactly the same argument shows that the eigenvalues of a positive semidefinite matrix are all positive or zero; the eigenvalues of a negative definite matrix are all strictly negative real numbers; and the eigenvalues of a negative semidefinite matrix are all ≤ 0 . Since the trace of a matrix is the sum of its eigenvalues, we see that *the trace of a positive definite matrix is a positive real number*. Since the determinant of a matrix is the product of its eigenvalues, we see that *the determinant of a positive definite matrix is a positive real number*. In particular, *positive definite matrices are always invertible*. In the negative definite case, the trace is always negative, but the determinant is negative iff n (the size of the matrix) is odd; the determinant is positive for n even. We will prove in §7.15 that a matrix $A \in M_n(\mathbb{C})$ is positive definite iff for all $k \in [n]$, the matrix consisting of the first k rows and columns of A has strictly positive determinant.

Every real number is either positive, negative, or zero. However, it is now apparent that the corresponding property is *not* true for Hermitian matrices. It is certainly possible for a Hermitian matrix to have both positive and negative real eigenvalues (assuming $n \geq 2$); consider, for example, diagonal matrices with real diagonal entries. Such matrices are neither positive definite nor negative definite. A Hermitian matrix with both positive and negative eigenvalues is sometimes called *indefinite*.

A matrix with all positive eigenvalues need not be Hermitian, and hence need not be positive definite; consider $A = \begin{bmatrix} 1 & 2 \\ 0 & 3 \end{bmatrix}$ for an example. We will see later that a *normal* matrix with all positive eigenvalues must be positive definite.

7.5 Unitary Matrices

Next we will develop a matrix analogue of the unit circle in \mathbb{C} . Recall that a complex number $z \in \mathbb{C}$ lies on the unit circle iff $|z| = 1$ iff $|z|^2 = 1$ iff $z\bar{z} = 1$ iff $z^{-1} = \bar{z}$. Accordingly, we define a matrix $U \in M_n(\mathbb{C})$ to be *unitary* iff $UU^* = I$. Note that a 1×1 matrix is unitary iff its sole entry has modulus 1. Set $u = \det(U) \in \mathbb{C}$; then $\det(U^*) = \bar{u}$. Taking determinants in the formula $UU^* = I$, we see that $|u|^2 = u\bar{u} = 1$; thus, *the determinant of a unitary matrix is a complex number of modulus 1*. In particular, $\det(U) \neq 0$, so U^{-1} exists. Multiplying $UU^* = I$ by U^{-1} , we see that $U^{-1} = U^*$ and hence $U^*U = I$. Reversing this argument, we see that the following three conditions on $U \in M_n(\mathbb{C})$ are equivalent:

$$UU^* = I; \quad U^*U = I; \quad U \text{ is invertible and } U^{-1} = U^*.$$

We could have taken any of these conditions as the definition of a unitary matrix.

If U is unitary, then U^* is also unitary, since $U^*(U^*)^* = U^*U = I$. If U is unitary, then U^{-1} is also unitary, since $U^{-1} = U^*$. If U is unitary, then the transpose U^T is also unitary, since $U^T(U^T)^* = U^T(U^*)^T = (U^*U)^T = I^T = I$. If U and V are unitary matrices of the same size, then UV is unitary, since $(UV)(UV)^* = UVV^*U^* = UIU^* = UU^* = I$. If U is a diagonal matrix such that each diagonal entry is a number $z_k \in \mathbb{C}$ of modulus 1, then U is unitary. For, $z_k^{-1} = \overline{z_k}$ implies $U^{-1} = U^*$ for such a matrix U . In particular, the identity matrix I is unitary. Let $U_n(\mathbb{C})$ denote the set of all unitary matrices in $M_n(\mathbb{C})$. The preceding remarks show that $U_n(\mathbb{C})$ is closed under identity, matrix product, and inverses. Therefore, the set $U_n(\mathbb{C})$ of unitary matrices forms a *subgroup* of the multiplicative group of invertible matrices $\text{GL}_n(\mathbb{C})$.

Unitary matrices preserve inner products; more precisely, if U is unitary and $v, w \in \mathbb{C}^n$, then $\langle Uv, Uw \rangle = \langle v, w \rangle$. This follows since $\langle Uv, Uw \rangle = (Uw)^*(Uv) = w^*(U^*U)v = w^*Iv = w^*v = \langle v, w \rangle$. Similarly, *unitary matrices preserve norms*, meaning that $\|Uv\| = \|v\|$ for all unitary U and all $v \in \mathbb{C}^n$. This follows since $\|Uv\|^2 = \langle Uv, Uv \rangle = \langle v, v \rangle = \|v\|^2$.

All eigenvalues of a unitary matrix are complex numbers of modulus 1. To see this, let $U \in M_n(\mathbb{C})$ be unitary with eigenvalue $\lambda \in \mathbb{C}$ and associated nonzero eigenvector v . We have $\|v\| = \|Uv\| = \|\lambda v\| = |\lambda| \cdot \|v\|$. Dividing by the nonzero scalar $\|v\|$ gives $|\lambda| = 1$, as needed. Taking the product of all the eigenvalues, we get another proof that $|\det(U)| = 1$ in \mathbb{C} .

A list of vectors $v_1, \dots, v_m \in \mathbb{C}^n$ is called *orthonormal* iff $\langle v_j, v_k \rangle = 0$ for all $j \neq k$ and $\langle v_j, v_j \rangle = 1$ for all j . Consider the condition $U^*U = I$, which can be used to define unitary matrices. Let the columns of U be $v_1, \dots, v_n \in \mathbb{C}^n$. The j, k -entry of the product U^*U is found by taking the product of row j of U^* and column k of U . But row j of U^* is v_j^* , and column k of U is v_k . We conclude that $(U^*U)(j, k) = v_j^*v_k = \langle v_k, v_j \rangle$ for all j, k . On the other hand, $I(j, k) = 1$ for $j = k$ and $I(j, k) = 0$ for $j \neq k$. Comparing to the definition of orthonormality, we see that *a matrix $U \in M_n(\mathbb{C})$ is unitary iff $U^*U = I$ iff the columns of U are orthonormal vectors in \mathbb{C}^n .* Applying similar reasoning to the condition $UU^* = I$ (or applying the result just stated to U^T), we see that *a matrix $U \in M_n(\mathbb{C})$ is unitary iff the rows of U are orthonormal vectors in \mathbb{C}^n .* Thus, the rows (resp. columns) of a unitary matrix U are mutually perpendicular unit vectors in \mathbb{C}^n , further strengthening the analogy between U and the unit circle in \mathbb{C} .

Suppose $U \in M_n(\mathbb{C})$ preserves inner products; i.e., $\langle Uv, Uw \rangle = \langle v, w \rangle$ for all $v, w \in \mathbb{C}^n$. Must U be unitary? The condition says that $w^*(U^*U)v = w^*v$ for all $v, w \in \mathbb{C}^n$. Choosing $v = e_j$ and $w = e_k$ (the standard basis vectors) and carrying out the matrix multiplication on the left side, we see that $(U^*U)(k, j)$ equals zero for $k \neq j$ and equals one for $k = j$. Thus, $U^*U = I$, and U is indeed unitary.

Suppose $U \in M_n(\mathbb{C})$ preserves norms; i.e., $\|Uv\| = \|v\|$ for all $v \in \mathbb{C}^n$. Must U be unitary? Here, the condition means that $\langle Uv, Uv \rangle = \langle v, v \rangle$ for all $v \in \mathbb{C}^n$. We will show that U must preserve inner products, hence is unitary. Fix $v, w \in \mathbb{C}^n$, and apply the assumed condition to the vectors v , w , $v + w$, and $v + iw$. We conclude that:

$$\langle Uv, Uv \rangle = \langle v, v \rangle; \quad \langle Uw, Uw \rangle = \langle w, w \rangle;$$

$$\langle U(v + w), U(v + w) \rangle = \langle v + w, v + w \rangle; \quad \langle U(v + iw), U(v + iw) \rangle = \langle v + iw, v + iw \rangle.$$

Now, for any complex scalars c, d and any vectors $v, w \in \mathbb{C}^n$, we have

$$\langle cv + dw, cv + dw \rangle = (cv + dw)^*(cv + dw) = \|cv\|^2 + \bar{d}c w^*v + \bar{c}d v^*w + \|dw\|^2.$$

Applying this formula to the preceding expressions, we get:

$$\begin{aligned}\langle v + w, v + w \rangle &= ||v||^2 + w^*v + v^*w + ||w||^2 \\ \langle Uv + Uw, Uv + Uw \rangle &= ||Uv||^2 + w^*U^*Uv + v^*U^*Uw + ||Uw||^2 \\ \langle v + iw, v + iw \rangle &= ||v||^2 - iw^*v + iv^*w + ||w||^2 \\ \langle Uv + iUw, Uv + iUw \rangle &= ||Uv||^2 - iw^*U^*Uv + iv^*U^*Uw + ||Uw||^2.\end{aligned}$$

For brevity, set $a = w^*v = \langle v, w \rangle$, $b = w^*(U^*U)v = \langle Uv, Uw \rangle$, and note that $v^*w = \bar{a}$ and $v^*(U^*U)w = \bar{b}$. Putting these values into the preceding formulas and recalling that $||Uv||^2 = ||v||^2$ and $||Uw||^2 = ||w||^2$, we get $a + \bar{a} = b + \bar{b}$ and $-i(a - \bar{a}) = -i(b - \bar{b})$. These equations say that $2\operatorname{Re}(a) = 2\operatorname{Re}(b)$ and $2\operatorname{Im}(a) = 2\operatorname{Im}(b)$; so the complex numbers a and b are equal. This means that $\langle Uv, Uw \rangle = \langle v, w \rangle$, so the argument in the last paragraph proves that U is unitary.

To summarize, all of the following conditions on a matrix $U \in M_n(\mathbb{C})$ are equivalent: U is unitary; $UU^* = I$; $U^*U = I$; U is invertible and $U^{-1} = U^*$; the rows of U are orthonormal in \mathbb{C}^n ; the columns of U are orthonormal in \mathbb{C}^n ; U^* is unitary; U^{-1} is unitary; U^T is unitary; U preserves inner products in \mathbb{C}^n ; U preserves norms in \mathbb{C}^n . Also, if U is unitary then the determinant and all eigenvalues of U have norm 1 in \mathbb{C} , but the converse does not hold in general. The failure of the converse can be demonstrated by considering upper-triangular matrices with numbers of modulus 1 on the main diagonal and at least one nonzero off-diagonal entry.

7.6 Unitary Similarity

Two square matrices $A, B \in M_n(\mathbb{C})$ are called *similar* iff there is an invertible matrix $S \in M_n(\mathbb{C})$ such that $B = S^{-1}AS$. We say A and B are *unitarily similar* iff there is a unitary matrix $U \in M_n(\mathbb{C})$ such that $B = U^{-1}AU = U^*AU$. One may check that similarity and unitary similarity are equivalence relations on $M_n(\mathbb{C})$.

Unitary similarity preserves the properties of being Hermitian, positive definite, or unitary. More precisely, if U is unitary and $B = U^*AU$, then A is Hermitian iff B is Hermitian; A is positive definite iff B is positive definite; and A is unitary iff B is unitary. For, if A is Hermitian, $B^* = (U^*AU)^* = U^*A^*U = U^*AU = B$. If A is positive definite and $v \in \mathbb{C}^n$ is nonzero, then $Uv \neq 0$ since U is invertible, hence $v^*Bv = v^*U^*AUv = (Uv)^*A(Uv)$ is a positive real number, so B is positive definite. Finally, if A is unitary, then $B = U^*AU$ is a product of unitary matrices and is therefore unitary. The converse statements follow from the symmetry of unitary similarity.

What is the geometric significance of similarity and unitary similarity? Suppose V is a finite-dimensional complex vector space (which we can take to be \mathbb{C}^n with no loss of generality), and suppose $T : V \rightarrow V$ is a linear operator on V . Suppose A is the matrix representing T relative to some ordered basis of V . If we switch to some other ordered basis of V , then the matrix representing T relative to the new basis is similar to A . The columns of the similarity matrix S give the coordinates of the new basis relative to the old basis. Conversely, any matrix similar to A is the matrix of T relative to some ordered basis of V . For more details, see Chapter 6.

Unitary similarity has an analogous geometric interpretation. This time we suppose V is a finite-dimensional complex inner product space (which we can take to be \mathbb{C}^n with the standard inner product, without loss of generality), and we assume T is a linear operator on V . Let A be the matrix of T relative to an *orthonormal* basis of V , which could be the

standard ordered basis (e_1, \dots, e_n) of \mathbb{C}^n . One checks that A is a length-preserving matrix (meaning $\|Ax\| = \|x\|$ for all $x \in \mathbb{C}^n$) iff T is a length-preserving map (meaning $\|Tv\| = \|v\|$ for all $v \in V$). If we switch to some other *orthonormal* basis of V , we obtain a new matrix B such that $B = U^{-1}AU$ for some invertible matrix U . As pointed out above, the columns of U are the coordinates of the new basis vectors relative to the old ones. Since both bases are orthonormal, we see that the columns of U are orthonormal in \mathbb{C}^n , and hence U is unitary. Thus, B is *unitarily similar* to A . Conversely, any matrix that is unitarily similar to A is the matrix of T relative to some *orthonormal* basis of V . So unitary similarity corresponds to *orthonormal* change of basis, just as ordinary similarity corresponds to *arbitrary* change of basis. For more details, see §6.21.

A matrix $A \in M_n(\mathbb{C})$ is called *unitarily diagonalizable* iff there exists a unitary matrix $U \in M_n(\mathbb{C})$ and a diagonal matrix $D \in M_n(\mathbb{C})$ such that $U^*AU = D$. A linear operator T on V is called *unitarily diagonalizable* iff the matrix of T relative to some orthonormal basis is diagonal iff the matrix of T relative to any orthonormal basis is unitarily diagonalizable. This means that there exists an orthonormal list of vectors (u_1, \dots, u_n) in V and scalars $c_1, \dots, c_n \in \mathbb{C}$ such that $T(u_k) = c_k u_k$ for all $k \in [n]$. In other words, *an operator T on V is unitarily diagonalizable iff the inner product space V has an orthonormal basis consisting of eigenvectors of T .*

In the case where $V = \mathbb{C}^n$ and A is the matrix of T relative to (e_1, \dots, e_n) , the matrix equation $U^*AU = D$ is equivalent to $AU = UD$ (since $U^* = U^{-1}$). Let the columns of U be $u_1, \dots, u_n \in \mathbb{C}^n$, and let the diagonal entries of D be c_1, \dots, c_n . The k 'th column of AU is Au_k , while the k 'th column of UD is $c_k u_k$. Thus, we get a matrix proof of the observation in the previous paragraph: a matrix $A \in M_n(\mathbb{C})$ is unitarily diagonalizable iff there exist n orthonormal vectors in \mathbb{C}^n that are eigenvectors of A . In this case, the unitary matrix achieving diagonalization has columns consisting of the orthonormal eigenvectors of A , and the resulting diagonal matrix has the eigenvalues of A on its main diagonal.

7.7 Unitary Triangularization

Not every matrix $A \in M_n(\mathbb{C})$ can be unitarily diagonalized. For example, the matrix $A = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$ has zero as its only eigenvalue. If A were similar (or unitarily similar) to a diagonal matrix, that matrix must have zeroes on its main diagonal, since similar matrices have the same eigenvalues. But then A would be similar to the zero matrix, forcing A itself to be zero. We will soon prove that *normality* is a necessary and sufficient condition for unitarily diagonalizability.

On the other hand, if we weaken the concept of unitary diagonalizability, we can obtain a result valid for all square matrices. Recall that a matrix $A \in M_n(\mathbb{C})$ is *triangulable* iff there is an invertible matrix $S \in M_n(\mathbb{C})$ such that $S^{-1}AS$ is an upper-triangular matrix. Call A *unitarily triangulable* iff there is a unitary $U \in M_n(\mathbb{C})$ such that $U^{-1}AU = U^*AU$ is an upper-triangular matrix. Geometrically, if $T : V \rightarrow V$ is the linear operator represented by the matrix A , then T is triangulable iff there exists an ordered basis (x_1, \dots, x_n) of V such that $T(x_i)$ lies in the span of (x_1, \dots, x_i) for all $i \in [n]$. T is unitarily triangulable iff there exists an orthonormal basis X of V with the stated property. Note, in this situation, that x_1 must be an eigenvector of T .

Schur proved the theorem that *every matrix $A \in M_n(\mathbb{C})$ can be unitarily triangularized*. Equivalently, for every linear operator T on an n -dimensional complex inner product space V , there is an orthonormal basis X of V such that $[T]_X$ (the matrix of T relative to X)

is upper-triangular. We will prove the result for linear operators by induction on $n \geq 1$. For $n = 1$, let x be any nonzero vector in V . Letting $X = (x/\|x\|)$, X is an orthonormal basis of V , and the 1×1 matrix $[T]_X$ is upper-triangular. Now fix $n > 1$, and assume the theorem is already known for all inner product spaces of dimension less than n . Since we are working over the algebraically closed field \mathbb{C} , we know that T has an eigenvalue $c_1 \in \mathbb{C}$ with associated eigenvector x_1 (take c_1 to be any root of the characteristic polynomial of T). Dividing x_1 by a scalar, we can ensure that $\|x_1\| = 1$. We claim $V = \mathbb{C}x_1 \oplus W$, where

$$W = x_1^\perp = \{w \in V : \langle w, x_1 \rangle = 0\}.$$

To see this, consider the map $S : V \rightarrow \mathbb{C}$ given by $S(v) = \langle v, x_1 \rangle$ for $v \in V$. S is a linear map with kernel W and image \mathbb{C} (since $S(x_1) = 1$), so the rank-nullity theorem shows that $n = \dim(V) = \dim(W) + \dim(\mathbb{C})$. As $\dim(\mathbb{C}) = 1$, we obtain $\dim(W) = n - 1$. On the other hand, given $dx_1 \in \mathbb{C}x_1 \cap W$ with $d \in \mathbb{C}$, we have $0 = S(dx_1) = \langle dx_1, x_1 \rangle = d\|x_1\|^2 = d$, so $\mathbb{C}x_1 \cap W = \{0\}$. This proves the claim.

It follows that every vector in V can be written uniquely in the form $ax_1 + w$, where $a \in \mathbb{C}$ and $w \in W$. Let $P : V \rightarrow W$ be the projection map defined by $P(ax_1 + w) = w$, and define a linear map $T' : W \rightarrow W$ by setting $T'(w) = P(T(w))$ for $w \in W$. By induction, we can find an orthonormal basis $X' = (x_2, \dots, x_n)$ of W such that $[T']_{X'}$ is upper-triangular. Let $X = (x_1, x_2, \dots, x_n)$, which is readily verified to be an orthonormal basis for V . We assert that

$$[T]_X = \left[\begin{array}{c|cccc} c_1 & * & * & \cdots & * \\ \hline 0 & & & & \\ \vdots & & [T']_{X'} & & \\ 0 & & & & \end{array} \right],$$

which is an upper-triangular matrix. To verify this, recall that the j 'th column of $[T]_X$ consists of the coordinates of $T(x_j)$ relative to the basis X . Now $T(x_1) = c_1 x_1$, so the first column of $[T]_X$ is $(c_1, 0, \dots, 0)$. For $j > 1$, write $T(x_j) = c_{1,j}x_1 + c_{2,j}x_2 + \cdots + c_{n,j}x_n$ with $c_{i,j} \in \mathbb{C}$. Then

$$T'(x_j) = P(T(x_j)) = c_{2,j}x_2 + \cdots + c_{n,j}x_n.$$

Thus, the j 'th column of $[T]_X$ consists of some scalar $c_{1,j}$ followed by the $(j-1)$ 'th column of $[T']_{X'}$. This completes the proof of Schur's theorem.

7.8 Simultaneous Triangularization

Suppose A and B are two matrices in $M_n(\mathbb{C})$. Schur's theorem guarantees that there are unitary matrices U and V such that U^*AU and V^*BV are upper-triangular matrices. However, there is no guarantee that the *same* matrix U will triangularize both A and B at the same time. If such a unitary matrix does exist, we will say that A and B can be *simultaneously* unitarily triangularized.

We will prove that *if $AB = BA$, then A and B can be simultaneously unitarily triangularized*. Equivalently, we shall prove: if T and U are two commuting operators on a complex inner product space V , then there exists an orthonormal basis X of V such that $[T]_X$ and $[U]_X$ are both upper-triangular matrices. Again we use induction on $n = \dim(V)$, the case $n = 1$ being immediate. Assume $n > 1$. If the required result is to be true, consideration of the first column of $[T]_X$ and $[U]_X$ shows that we need to find a unit vector $x_1 \in \mathbb{C}^n$ that is simultaneously an eigenvector for T and an eigenvector for U . To do this,

fix an eigenvalue c_1 of T . Let $Z = \{v \in V : T(v) = c_1 v\}$, which is a nonzero subspace of V (Z is the eigenspace associated with the eigenvalue c_1). We claim that U maps Z into Z . For if $v \in Z$, commutativity of T and U gives

$$T(U(v)) = U(T(v)) = U(c_1 v) = c_1(U(v)),$$

so that $U(v) \in Z$. It follows that the restriction $U|Z : Z \rightarrow Z$ is a linear operator on Z . Thus, U has an eigenvector in Z (with associated eigenvalue d_1 , say). Normalizing if necessary, let $x_1 \in Z$ be a unit eigenvector of U . Since $x_1 \in Z$, x_1 is also an eigenvector for T .

Now we continue as in the previous proof. Let $W = x_1^\perp$, so that $V = \mathbb{C}x_1 \oplus W$. Consider the two linear operators on W defined by $T' = (P \circ T)|W$ and $U' = (P \circ U)|W$, where $P : \mathbb{C}x_1 \oplus W \rightarrow W$ is the projection $P(ax_1 + w) = w$. We want to apply our induction hypothesis to the $(n-1)$ -dimensional space W . Before doing so, we must check that the new operators T' and U' on W still commute. To verify this, first note that $P \circ T \circ P = P \circ T$. For, given any $v = ax_1 + w$ in V (where $a \in \mathbb{C}$ and $w \in W$),

$$\begin{aligned} P(T(v)) &= P(T(ax_1 + w)) = P(aT(x_1) + T(w)) = \\ &\quad P(ac_1 x_1) + P(T(w)) = P(T(w)) = P(T(P(v))). \end{aligned} \quad (7.2)$$

Similarly, $P \circ U \circ P = P \circ U$. It follows that

$$\begin{aligned} T' \circ U' &= (P \circ T \circ P \circ U)|W = (P \circ T \circ U)|W \\ &= (P \circ U \circ T)|W = (P \circ U \circ P \circ T)|W = U' \circ T'. \end{aligned}$$

The induction hypothesis can now be applied to obtain an orthonormal basis $X' = (x_2, \dots, x_n)$ of W such that $[T']_{X'}$ and $[U']_{X'}$ are both upper-triangular. Letting $X = (x_1, x_2, \dots, x_n)$, we see as in the previous proof that X is an orthonormal basis of V such that $[T]_X$ and $[U]_X$ are both upper-triangular. This completes the induction.

This result extends to an arbitrary (possibly infinite) family of commuting matrices or operators. Specifically, suppose V is an n -dimensional complex inner product space and $\{T_i : i \in I\}$ is an indexed family of operators on V such that $T_i \circ T_j = T_j \circ T_i$ for all $i, j \in I$. Then there exists an orthonormal basis X of V such that $[T_i]_X$ is upper-triangular for all $i \in I$. We can prove this result by imitating the previous induction proof. The key point is that we must locate a unit vector in V that is simultaneously an eigenvector for *every* operator T_i . To do this, consider a nonzero subspace Z of V of minimum positive dimension such that $T_i[Z] \subseteq Z$ for all $i \in I$. Such subspaces do exist — at worst, $Z = V$. Let T be a fixed operator from the family $\{T_i : i \in I\}$. Let Z' be a nonzero eigenspace for $T|Z$. The argument used before shows that $T_i[Z'] \subseteq Z'$ for all $i \in I$, since each T_i commutes with T . By minimality of $\dim(Z)$, we see that $Z' = Z$. Thus, Z itself is an eigenspace for T , so every nonzero $z \in Z$ is an eigenvector for T . But T was arbitrary, so we see that every nonzero $z \in Z$ is an eigenvector for every T_i in the given family. Choose any such z and normalize it to get a common unit eigenvector x_1 for all the T_i 's. The rest of the proof (see the previous paragraph) now goes through verbatim. In particular, the family of maps $\{T'_i = (P \circ T_i)|W : i \in I\}$ is still a commuting family (by the same calculation as before), so the induction hypothesis does apply to give us the basis X' as above.

7.9 Normal Matrices and Unitary Diagonalization

Recall that a matrix A is unitarily diagonalizable iff there exists a unitary matrix U such that U^*AU is a diagonal matrix. Equivalently, an operator T on a complex inner product space V is unitarily diagonalizable iff there exists some orthonormal basis X of V such that $[T]_X$ is diagonal, i.e., such that the elements of X are all eigenvectors of T . We want to find an easily checked necessary and sufficient condition on A (or T) that characterizes when A can be unitarily diagonalized.

Recall that a matrix $A \in M_n(\mathbb{C})$ is *normal* iff $AA^* = A^*A$. This condition arose in our consideration of the Hermitian decomposition $A = X + iY$. We saw in §7.3 that $XY = YX$ iff A is normal. We can see quickly that normality is a *necessary* condition for unitary diagonalizability. For, suppose U is a unitary matrix and D is a diagonal matrix such that $D = U^*AU$. Then $D^* = U^*A^*U$ is also diagonal, and hence $DD^* = D^*D$. But then

$$AA^* = (UDU^*)(UD^*U^*) = U(DD^*)U^* = U(D^*D)U^* = (UD^*U^*)(UDU^*) = A^*A,$$

so A is normal.

We shall show that normality is also *sufficient* for unitary diagonalizability, so that a matrix $A \in M_n(\mathbb{C})$ is unitarily diagonalizable iff A is normal. This fundamental result is often called the *spectral theorem for normal matrices*. To prove sufficiency, suppose $A \in M_n(\mathbb{C})$ is a normal matrix. Then $AA^* = A^*A$ by definition, so A and A^* can be simultaneously unitarily triangularized. Choose a unitary matrix U such that $C = U^*AU$ and $D = U^*A^*U$ are both upper-triangular. On one hand, C is an upper-triangular matrix. On the other hand, $C = U^*AU = (U^*A^*U)^* = D^*$ is the conjugate-transpose of an upper-triangular matrix, so C is a lower-triangular matrix. This means that C must be a diagonal matrix, so we have unitarily diagonalized A .

What are some examples of normal matrices? First, *Hermitian matrices are normal*, since $A = A^*$ implies $AA^* = A^2 = A^*A$. In particular, positive and negative definite (or semidefinite) matrices are normal. Second, *unitary matrices are normal*, since for U unitary, $UU^* = I = U^*U$. Third, skew-Hermitian matrices are normal, since $A^* = -A$ implies $AA^* = -A^2 = A^*A$. Fourth, if A is normal and U is unitary, then U^*AU is normal, as the reader may verify. One can readily manufacture normal matrices that are neither Hermitian, skew-Hermitian, nor unitary. For example, consider $A = \begin{bmatrix} 1+i & 0 \\ 0 & 4 \end{bmatrix}$ or any matrix unitarily similar to A .

Applying the spectral theorem to these examples, we see that *Hermitian matrices, positive definite matrices, and unitary matrices can all be unitarily diagonalized*. In particular, the linear operators corresponding to such matrices are guaranteed to have an associated orthonormal basis of eigenvectors. The nature of the eigenvalues of a normal matrix can be used to characterize the properties of being Hermitian, positive definite, or unitary. More precisely, if $A \in M_n(\mathbb{C})$ is normal, then:

- A is Hermitian iff all eigenvalues of A are real;
- A is positive definite iff all eigenvalues of A are strictly positive real numbers;
- A is unitary iff all eigenvalues of A are complex numbers of modulus 1;
- A is invertible iff all eigenvalues of A are nonzero.

In each case, the forward implication holds with no restriction on A (as seen earlier). It is now straightforward to prove both directions of each implication using the spectral theorem.

Given a normal matrix A , write $D = U^*AU$ where U is unitary and D is a diagonal matrix with the eigenvalues of A on its diagonal. Unitary similarity preserves the properties of being Hermitian, positive definite, unitary, or invertible, so A has one of these properties iff D does. But it is immediate from the definitions that a diagonal matrix D is Hermitian iff all diagonal entries of D are real; D is positive definite iff all diagonal entries of D are strictly positive; D is unitary iff all diagonal entries of D have norm 1 in \mathbb{C} ; and D is invertible iff all diagonal entries of D are nonzero.

Here are a few more consequences of the spectral theorem. We know that every positive real number has a positive square root. This suggests the following matrix result: if $A \in M_n(\mathbb{C})$ is positive definite, there exists a positive definite matrix B such that $B^2 = A$. To prove this, write $D = U^*AU$ with D diagonal and U unitary. Each diagonal entry of D is a positive real number; let E be the diagonal matrix obtained by replacing each such entry by its positive square root. Then $E^2 = D$. Setting $B = UEU^*$, a routine computation confirms that $B^2 = A$. Moreover, since E is positive definite, so is B . A similar argument shows that, for example, every Hermitian matrix A has a Hermitian “cube root,” which is a Hermitian matrix C such that $C^3 = A$.

7.10 Polynomials and Commuting Matrices

One sometimes needs to know which matrices in $M_n(\mathbb{C})$ commute with a given matrix A . For example, we have already seen that commutativity of A and B is a sufficient condition for A and B to be simultaneously triangulable. In general, it is difficult to characterize the matrices that commute with A in any simple way. However, we can describe a collection of matrices that always commute with A . Consider a polynomial $p \in \mathbb{C}[x]$ with complex coefficients, say $p = c_0 + c_1x + c_2x^2 + \cdots + c_dx^d$ where all $c_i \in \mathbb{C}$. We can replace x by A in this polynomial to obtain a matrix $p(A) = c_0I + c_1A + c_2A^2 + \cdots + c_dA^d \in M_n(\mathbb{C})$. We claim that A commutes with $p(A)$. For,

$$Ap(A) = c_0A + c_1A^2 + c_2A^3 + \cdots + c_dA^{d+1} = p(A)A.$$

Thus, any matrix A commutes with every polynomial in A .

More generally, suppose C is any matrix that commutes with a given matrix A . We show by induction that C commutes with A^k for all integers $k \geq 0$. This holds for $k = 0$ and $k = 1$, since C commutes with I and A . Assuming the result holds for some k , note that

$$CA^{k+1} = (CA^k)A = (A^kC)A = A^k(CA) = A^k(AC) = A^{k+1}C.$$

Next, observe that if C commutes with U and V , then C commutes with $U + V$. If C commutes with U and r is a scalar, then C commutes with rU . It follows that if C commutes with matrices U_0, \dots, U_m , then C commutes with any linear combination of U_0, \dots, U_m . Taking $U_k = A^k$, we deduce that *if C commutes with A , then C commutes with $p(A)$ for every polynomial p* . For example, if A is invertible, then $C = A^{-1}$ commutes with A (since $AA^{-1} = I = A^{-1}A$), so A^{-1} commutes with every polynomial in A .

We have just shown that for all matrices $A, C \in M_n(\mathbb{C})$, if C commutes with A , then $p(A)$ commutes with C for every polynomial p . We can apply this result with A replaced by C and C replaced by $p(A)$ to conclude: *if $AC = CA$, then $p(A)$ commutes with $q(C)$ for all polynomials $p, q \in \mathbb{C}[x]$* .

Here are a few more facts that can be proved by similar methods. First, suppose $B = U^*AU$ for some unitary matrix U . Then $p(B) = U^*p(A)U$ and $p(A) = Up(B)U^*$ for all

polynomials $p \in \mathbb{C}[x]$. As above, we first prove this for the particular polynomials $p = x^k$ by induction on $k \geq 0$, and then deduce the result for general p by linearity. Second, suppose D is a diagonal matrix with diagonal entries d_1, \dots, d_n . The same proof technique shows that $p(D)$ is a diagonal matrix with diagonal entries $p(d_1), \dots, p(d_n)$.

For our next theorem, we will show that A is normal iff $A^* = p(A)$ for some polynomial $p \in \mathbb{C}[x]$. One direction is immediate: if $A^* = p(A)$, then A^* commutes with A , so A is normal. To prove the other direction, assume A is normal. Choose a unitary U such that $D = U^*AU$ is diagonal; let (d_1, \dots, d_n) be the diagonal entries of D . We know that $U^*A^*U = D^*$, where D^* has diagonal entries $(\overline{d_1}, \dots, \overline{d_n})$. Suppose we can find a polynomial p such that $p(D) = D^*$. A routine calculation shows that $p(A) = p(UDU^*) = Up(D)U^* = UD^*U^* = A^*$, so the proof will be complete once p is found.

To construct p , we use the following algebraic fact called *Lagrange's interpolation formula* (see §3.18). If x_1, \dots, x_m are m distinct numbers in \mathbb{C} , and y_1, \dots, y_m are arbitrary elements of \mathbb{C} , then there exists a polynomial $p \in \mathbb{C}[x]$ such that $p(x_i) = y_i$ for all i . Indeed, one readily checks that

$$p(x) = \sum_{j=1}^m y_j \frac{\prod_{i:i \neq j}(x - x_i)}{\prod_{i:i \neq j}(x_j - x_i)}$$

is a polynomial with the required values. To apply this theorem in the current situation, note that $d_i = d_j$ iff $\overline{d_i} = \overline{d_j}$. We can take the x_i 's to be the *distinct* complex numbers that occur on the main diagonal of D , and we can take $y_i = \overline{x_i}$ for each i . Choosing p as above, it then follows that $p(d_i) = \overline{d_i}$ for all i . Accordingly, $p(D) = D^*$, and the proof is complete.

A similar proof shows that *if A is normal and invertible, then A^{-1} is a polynomial in A* . Write $D = U^*AU$ as above, and note that D^{-1} is the diagonal matrix with diagonal entries $d_1^{-1}, \dots, d_n^{-1}$. Using Lagrange interpolation as before, we can find a polynomial $q \in \mathbb{C}[x]$ such that $q(d_i) = d_i^{-1}$ for all i . Then $q(D) = D^{-1}$, and hence $q(A) = A^{-1}$.

7.11 Simultaneous Unitary Diagonalization

We say that a family of matrices $\{A_i : i \in I\} \subseteq M_n(\mathbb{C})$ can be *simultaneously unitarily diagonalized* iff there exists a unitary matrix U (independent of i) such that $D_i = U^*A_iU$ is diagonal for all $i \in I$. We are now ready to derive a necessary and sufficient condition for a family of matrices to be simultaneously unitarily diagonalizable.

It is certainly necessary that each individual A_i be normal (unitarily diagonalizable). Another quickly established necessary condition is that $A_iA_j = A_jA_i$ for all $i, j \in I$, i.e., the given family must be a *commuting* family of matrices. To see this, suppose $D_i = U^*A_iU$ is diagonal for all $i \in I$. Diagonal matrices commute, so $D_iD_j = D_jD_i$ for all $i, j \in I$. Therefore,

$$A_iA_j = (UD_iU^*)(UD_jU^*) = U(D_iD_j)U^* = U(D_jD_i)U^* = (UD_jU^*)(UD_iU^*) = A_jA_i$$

for all $i, j \in I$.

We will now prove that the necessary conditions in the preceding paragraph are also sufficient. In other words, *a family $\{A_i : i \in I\}$ of matrices in $M_n(\mathbb{C})$ is simultaneously unitarily diagonalizable iff every A_i is normal and $A_iA_j = A_jA_i$ for all $i, j \in I$* . To prove sufficiency, we recycle some ideas from our earlier proof that a normal matrix is unitarily diagonalizable. Recall that the key idea in that proof was to simultaneously unitarily upper-triangularize A and A^* , and then observe that the resulting triangular matrices must actually be diagonal. To use this idea here, we need to invoke our previous theorem about

simultaneous upper-triangularization of a family of matrices (§7.8). Consider the family of matrices $\mathcal{F} = \{A_i : i \in I\} \cup \{A_i^* : i \in I\}$. Each A_i is normal, so $A_i^* = p_i(A_i)$ for certain polynomials $p_i \in \mathbb{C}[x]$. To use our theorem, we need to know that \mathcal{F} is a commuting family. For any $i, j \in I$, A_i commutes with A_j by hypothesis. So, by our results on polynomials, A_i commutes with $p_j(A_j) = A_j^*$; $p_i(A_i) = A_i^*$ commutes with A_j ; and $p_i(A_i) = A_i^*$ commutes with $p_j(A_j) = A_j^*$. This covers all possible pairs of matrices in \mathcal{F} . So our theorem does apply, and we know there is a unitary matrix U such that $D_i = U^* A_i U$ and $E_i = U^* A_i^* U$ are upper-triangular matrices for all $i \in I$. We now finish exactly as before: D_i is both upper-triangular and lower-triangular, since $D_i = E_i^*$; so D_i is a diagonal matrix for all $i \in I$. Thus the family $\{A_i : i \in I\}$ has been simultaneously unitarily diagonalized. (In fact, we have even diagonalized the larger family \mathcal{F} .)

As an illustration of this theorem, let us prove a sharper version of an earlier result concerning matrix square roots. Let $A \in M_n(\mathbb{C})$ be a positive definite matrix. We will show that there exists a *unique* positive definite matrix B such that $B^2 = A$. Recall how existence of B follows from the spectral theorem: first, we find a unitary U such that $D = U^* A U$ is diagonal with positive real diagonal entries d_1, \dots, d_n . Then we let E be diagonal with diagonal entries $\sqrt{d_1}, \dots, \sqrt{d_n}$, and put $B = U E U^*$. By Lagrange interpolation, choose a polynomial $p \in \mathbb{C}[x]$ such that $p(d_i) = \sqrt{d_i}$ for all i . It follows that $p(D) = E$ and hence $p(A) = B$. Now, suppose C is any positive definite matrix such that $C^2 = A$. Since $B = p(C^2)$ is a polynomial in C , B and C commute. Both B and C are normal, so we can simultaneously unitarily diagonalize B and C . Let V be a unitary matrix such that $D_1 = V^* B V$ and $D_2 = V^* C V$ are both diagonal. D_1 and D_2 are positive definite (since B and C are), so all diagonal entries of D_1 and D_2 are positive. Furthermore, $D_1^2 = V^* A V = D_2^2$, so that $D_1(k, k)^2 = D_2(k, k)^2 \in \mathbb{R}$ for all $k \in [n]$. Since every positive real number has a unique positive square root, it follows that $D_1(k, k) = D_2(k, k)$ for all $k \in [n]$. So $D_1 = D_2$, which implies $B = V D_1 V^* = V D_2 V^* = C$. A similar argument shows that for all integers $k \geq 1$, every positive definite (resp. semidefinite) matrix in $M_n(\mathbb{C})$ has a unique positive definite (resp. semidefinite) $2k$ 'th root, and every Hermitian matrix has a unique Hermitian $(2k-1)$ 'th root. Hence, we may safely use the notation \sqrt{A} to denote the unique positive definite square root of a positive definite matrix A , $\sqrt[3]{B}$ to denote the unique Hermitian cube root of a Hermitian matrix B , etc.

7.12 Polar Decomposition: Invertible Case

For every nonzero complex number z , there exist $r, u \in \mathbb{C}$ such that r is a positive real number, $|u| = 1$, and $z = ru = ur$. We call this the *polar decomposition* of z . Analogy suggests the following result: for every invertible normal matrix $A \in M_n(\mathbb{C})$, there exist $R, U \in M_n(\mathbb{C})$ such that R is positive definite, U is unitary, and $A = RU = UR$. The techniques developed so far lead to a quick proof, as follows. By normality of A , find a unitary V and a diagonal D such that $D = V^* A V$. Each diagonal entry $d_i = D(i, i)$ is a nonzero complex number (since A is invertible). Let $d_i = r_i u_i$ be the polar decomposition of d_i . Let R' and U' be diagonal matrices with diagonal entries r_i and u_i , so R' is positive definite and U' is unitary. Observe that $D = R' U' = U' R'$. Putting $R = VR'V^*$ and $U = VU'V^*$, we then have $A = RU = UR$ where R is positive definite and U is unitary.

More generally, if A is *any* invertible matrix in $M_n(\mathbb{C})$ (not necessarily normal), there exist *unique* matrices $R, U \in M_n(\mathbb{C})$ such that R is positive definite, U is unitary, and $A = RU$. This factorization is called the *polar decomposition* of A . To motivate the proof, return to the polar decomposition $z = ru$ of a nonzero complex number z . To find r given

z , we could calculate $z\bar{z} = |z|^2 = r^2|u|^2 = r^2$ and then take square roots to find r . Then u must be $r^{-1}z$. Given the invertible matrix A , we are therefore led to consider the matrix AA^* . This matrix is positive definite, since for nonzero $v \in \mathbb{C}^n$, $v^*AA^*v = ||A^*v||^2$ is a positive real number. Let R be the unique positive definite square root of AA^* , and let $U = R^{-1}A$. Then $A = RU$. Since $R^* = R$ and $RR = AA^*$, the fact that U is unitary follows from the calculation

$$UU^* = R^{-1}AA^*(R^{-1})^* = R^{-1}(RR)(R^*)^{-1} = (R^{-1}R)(RR^{-1}) = I.$$

To see that R and U are uniquely determined by A , suppose $A = SW$ where S is positive definite and W is unitary. Then $AA^* = SWW^*S^* = SIS = S^2$. So $S = R$ since the positive definite matrix AA^* has a *unique* positive definite square root. It follows that $W = S^{-1}A = R^{-1}A = U$.

In the polar decomposition $A = RU$, R and U need not commute. In fact, $RU = UR$ holds iff A is a normal matrix (compare this to the corresponding fact for the Cartesian decomposition of A , which motivated the definition of normality). If A is normal, $RU = UR$ follows from the unitary diagonalization proof given in the first paragraph of this section. Conversely, suppose R and U commute. Since R and U are both normal, it suffices to prove that *the product of two commuting normal matrices is also normal*. Suppose B and C are normal commuting matrices. By normality, B^* is a polynomial in B and C^* is a polynomial in C , so the four matrices B, B^*, C, C^* all commute with one another. Then

$$(BC)(BC)^* = BCC^*B^* = C^*B^*BC = (BC)^*(BC),$$

so BC is normal.

For A invertible but not normal, we can obtain a “dual” polar decomposition $A = U_1R_1$, where R_1 is positive definite and U_1 is unitary. It suffices to take the conjugate-transpose of the polar decomposition $A^* = R_2U_2$, which gives $A = U_2^*R_2^*$ where U_2^* is unitary and $R_2^* = R_2$ is positive definite. As above, U_1 and R_1 are uniquely determined by A .

7.13 Polar Decomposition: General Case

The complex number 0 can be written in polar form as $0 = ru = ur$, where $r = 0$ and u is *any* complex number of modulus 1. Analogy suggests the following result: for any non-invertible matrix $A \in M_n(\mathbb{C})$, there exist a positive semidefinite R and a unitary U such that $A = UR$, where R (but not U) is uniquely determined by A . Finding R and proving its uniqueness is not difficult, given what we did earlier for invertible A 's: if $A = UR$ is to hold, we must have $A^*A = R^*U^*UR = R^2$, forcing R to be the unique positive semidefinite square root of the matrix A^*A (which is readily seen to be positive semidefinite). However, finding U becomes trickier, since neither A nor R is invertible in the present situation.

To proceed, we need the observation that R and A have the same null space. To see this, fix $v \in \mathbb{C}^n$, and compute (using $R^2 = A^*A$ and $R = R^*$)

$$Rv = 0 \Leftrightarrow ||Rv||^2 = 0 \Leftrightarrow v^*R^*Rv = 0 \Leftrightarrow v^*A^*Av = 0 \Leftrightarrow ||Av||^2 = 0 \Leftrightarrow Av = 0.$$

By the rank-nullity theorem, the image (range) of R and the image of A have the same dimension. Let $V \subseteq \mathbb{C}^n$ be the image of R , and let $W \subseteq \mathbb{C}^n$ be the image of A . We define a map $\phi : V \rightarrow W$ as follows. Given $z \in V$, write $z = Rx$ for some $x \in \mathbb{C}^n$, and set $\phi(z) = Ax$. We must check that ϕ is well-defined, i.e., that the value $\phi(z)$ does not depend

on the choice of x such that $Rx = z$. If x' is another vector such that $Rx' = z$, then $R(x - x') = 0$, hence $A(x - x') = 0$, so $Ax = Ax'$. Thus, ϕ is well-defined, and this map satisfies $\phi(Rx) = Ax$ for all $x \in \mathbb{C}^n$. It is now quickly checked that ϕ is a linear map. Furthermore, ϕ is length-preserving, since

$$\|\phi(Rx)\|^2 = \|Ax\|^2 = x^* A^* Ax = x^* R^2 x = x^* R^* Rx = \|Rx\|^2$$

for all $x \in \mathbb{C}^n$. We can write $\mathbb{C}^n = V \oplus V^\perp = W \oplus W^\perp$ (see Exercise 24). Extend ϕ to a linear operator u on \mathbb{C}^n by sending an orthonormal basis of V^\perp to an orthonormal basis of W^\perp and extending by linearity; the operator u is readily seen to be length-preserving. Finally, let U be the matrix of u relative to the standard ordered basis of \mathbb{C}^n , so that $u(x) = Ux$ (matrix-vector product) for all $x \in \mathbb{C}^n$. Since u is a length-preserving map, U is a length-preserving (unitary) matrix. Furthermore, $URx = u(Rx) = \phi(Rx) = Ax$ for all $x \in \mathbb{C}^n$. Thus, $A = UR$, and the proof is complete. Moreover, the proof reveals the extent of the non-uniqueness of U : the action of U is forced on the image V of $R = \sqrt{A^* A}$, but U can act as an arbitrary length-preserving map from V^\perp to W^\perp .

Reasoning as in the invertible case, $UR = RU$ holds iff A is normal. By forming the conjugate-transpose of the polar decomposition $A^* = U_2 R_2$, we see that every $A \in M_n(\mathbb{C})$ can be written as a product $R_1 U_1$, where $R_1 = \sqrt{AA^*}$ is positive semidefinite and U_1 is unitary.

The polar decomposition $A = UR$ is often written in the following equivalent form. Unitarily diagonalize R , say $D = Q^* R Q$ for some unitary matrix Q . Then $R = Q D Q^*$, and $P = UQ$ is a unitary matrix. It follows that *any matrix $A \in M_n(\mathbb{C})$ can be factored in the form*

$$A = PDQ^*,$$

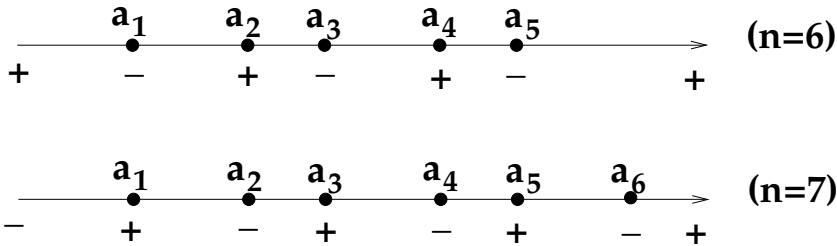
where Q and P are unitary, and D is a diagonal matrix with nonnegative real entries. This factorization is called the *singular value decomposition* of A . The entries on the diagonal of D are the eigenvalues of R (occurring with the correct multiplicities), and R is uniquely determined by A . It follows that the diagonal entries of D (disregarding order) are uniquely determined by the matrix A ; these entries are called the *singular values* of A . By passing from matrices to the linear maps they represent, we obtain the following geometric interpretation of the matrix factorization $A = PDQ^*$. Every linear transformation on \mathbb{C}^n is the composition of an isometry (length-preserving linear map), followed by a rescaling of the coordinate axes by nonnegative real numbers, followed by another isometry. The multiset of rescaling factors is uniquely determined by the linear map. We give another proof of the singular value decomposition, extended to rectangular matrices, in §9.12.

7.14 Interlacing Eigenvalues for Hermitian Matrices

Given a Hermitian matrix $B \in M_n(\mathbb{C})$ with $n > 1$, let $A \in M_{n-1}(\mathbb{C})$ be the submatrix of B obtained by deleting the last row and column of B . Evidently A is also Hermitian. We know all eigenvalues of A and B are real numbers. Let the eigenvalues of B be listed in increasing order as $b_1 \leq b_2 \leq \dots \leq b_n$, and let the eigenvalues of A be $a_1 \leq a_2 \leq \dots \leq a_{n-1}$. Our goal in this section is to prove the *interlacing theorem*

$$b_1 \leq a_1 \leq b_2 \leq a_2 \leq b_3 \leq a_3 \leq \dots \leq b_{n-1} \leq a_{n-1} \leq b_n,$$

which says that eigenvalues of B and A occur alternately as we scan the real number line from left to right. We proceed in steps, proving the result for successively more general matrices B .

**FIGURE 7.1**

Signs of $\chi_B(t)$ at the Eigenvalues of A .

Step 1. We prove the theorem for all matrices B of the form

$$B = \left[\begin{array}{cccc|c} a_1 & & & & c_1 \\ & a_2 & & & c_2 \\ & & \ddots & & \vdots \\ & & & a_{n-1} & c_{n-1} \\ \hline \bar{c}_1 & \bar{c}_2 & \dots & \bar{c}_{n-1} & a_n \end{array} \right], \quad (7.3)$$

where A is diagonal, a_1, \dots, a_{n-1} are distinct, $a_n \in \mathbb{R}$, and every $c_k \in \mathbb{C}$ is nonzero. Using the definition of determinants in §5.2, one first computes that any matrix of the form (7.3) has characteristic polynomial

$$\chi_B(t) = \det(tI_n - B) = \prod_{k=1}^n (t - a_k) - \sum_{k=1}^{n-1} (t - a_1) \cdots (t - a_{k-1}) |c_k|^2 (t - a_{k+1}) \cdots (t - a_{n-1}) \quad (7.4)$$

(see Exercise 57). We also know $\chi_B(t) = \prod_{k=1}^n (t - b_k)$ is a monic polynomial of degree n in $\mathbb{R}[t]$ whose roots are the eigenvalues of B (§5.15).

To proceed, we evaluate the polynomial $\chi_B(t)$ at $t = a_j$, where $1 \leq j \leq n-1$. Taking $t = a_j$ in (7.4), the first product becomes zero, and all summands for $k \neq j$ also become zero because of the factor $(t - a_j)$. Thus,

$$\chi_B(a_j) = -(a_j - a_1) \cdots (a_j - a_{j-1}) |c_j|^2 (a_j - a_{j+1}) \cdots (a_j - a_{n-1}).$$

Since the a_k 's are distinct and occur in increasing order, and since $|c_j|^2 \neq 0$, we see that $\chi_B(a_j)$ is negative for $j = n-1, n-3, n-5, \dots$, and $\chi_B(a_j)$ is positive for $j = n-2, n-4, n-6, \dots$. On the other hand, since $\chi_B(t)$ is monic of degree n , $\lim_{x \rightarrow \infty} \chi_B(x) = +\infty$, whereas $\lim_{x \rightarrow -\infty} \chi_B(x)$ is $+\infty$ for n even, and $-\infty$ for n odd. We illustrate this situation in Figure 7.1 in the cases $n = 6$ and $n = 7$. Since polynomial functions are continuous, we can apply the intermediate value theorem from calculus to conclude that each of the n intervals $(-\infty, a_1), (a_k, a_{k+1})$ for $1 \leq k < n-1$, and (a_{n-1}, ∞) contains at least one root of $\chi_B(t)$. Since $\chi_B(t)$ has exactly n real zeroes, we see that each interval contains exactly one eigenvalue of B . Then $b_1 < a_1 < b_2 < a_2 < \cdots < b_{n-1} < a_{n-1} < b_n$, as needed.

Step 2. We generalize step 1 by dropping the assumption that a_1, \dots, a_{n-1} are distinct. We show that for each multiple eigenvalue of A of multiplicity $s+1$, say $a_j = a_{j+1} = \cdots = a_{j+s}$, a_j is an eigenvalue of B of multiplicity s . Inspection of the right side of (7.4) shows that the polynomial $\chi_B(t)$ is divisible by $(t - a_j)^s$. On the other hand, dividing $\chi_B(t)$ by

$(t - a_j)^s$ and then setting $t = a_j$, we obtain

$$\begin{aligned} & - \sum_{k=j}^{j+s} (a_k - a_1) \cdots (a_k - a_{j-1}) |c_k|^2 (a_k - a_{j+s+1}) \cdots (a_k - a_{n-1}) \\ & = -(a_j - a_1) \cdots (a_j - a_{j-1}) (a_j - a_{j+s+1}) \cdots (a_j - a_{n-1}) (|c_j|^2 + \cdots + |c_{j+s}|^2) \neq 0. \end{aligned}$$

Thus, a_j is a root of $\chi_B(t)$ of multiplicity s . The proof of step 2 can now be completed by a sign-change analysis similar to step 1 (see Exercise 58).

Step 3. We generalize step 2 by dropping the assumption that every c_k be nonzero. Proceed by induction on $n \geq 2$. If $n = 2$ and $c_1 = 0$, the eigenvalues of B are a_1 and a_2 (in some order), and the needed interlacing inequalities immediately follow. If $n > 2$ and some $c_k = 0$, let B' (resp. A') be obtained from B (resp. A) by deleting row k and column k . Since $c_k = 0$, we see from (7.3) that $\chi_B(t) = (t - a_k)\chi_{B'}(t)$ and $\chi_A(t) = (t - a_k)\chi_{A'}(t)$. By induction, the eigenvalues of B' and A' satisfy the required interlacing property. To obtain the eigenvalues of B and A , we insert one more copy of a_k into both lists of eigenvalues. One quickly verifies that the interlacing property still holds for the new lists.

Step 4. We prove the result for general Hermitian $B \in M_n(\mathbb{C})$. As A is unitarily diagonalizable, there is a unitary matrix $U \in M_{n-1}(\mathbb{C})$ such that U^*AU is diagonal with entries a_1, \dots, a_{n-1} on the diagonal. Add a zero row at the bottom of U , and then add the column e_n on the right end to get a unitary matrix $V \in M_n(\mathbb{C})$. The matrix V^*BV has the same eigenvalues as B and is of the form (7.3). By step 3, the required interlacing property holds.

7.15 Determinant Criterion for Positive Definite Matrices

Given $A \in M_n(\mathbb{C})$, one can quickly decide whether A is Hermitian ($A^* = A$), unitary ($A^*A = I$), or normal ($A^*A = AA^*$) via matrix computations that do not require knowing the eigenvalues of A . On the other hand, the definition of a positive definite matrix ($v^*Av \in \mathbb{R}^+$ for all nonzero $v \in \mathbb{C}^n$) cannot be checked so easily. In this section, we prove a determinant condition on A that is necessary and sufficient for A to be positive definite. While computing determinants of large matrices is time-consuming, it may be even harder to find all the eigenvalues of A .

Assume $A \in M_n(\mathbb{C})$ is positive definite. We have already remarked that the trace and determinant of A must be positive real numbers (§7.4). More generally, for $1 \leq k \leq n$, let $A[k]$ denote the matrix in $M_k(\mathbb{C})$ consisting of the first k rows and columns of A . We claim $A[k]$ is positive definite, so that $\det(A[k]) \in \mathbb{R}^+$. To see this, fix a nonzero $v \in \mathbb{C}^k$. Let $w \in \mathbb{C}^n$ be v with $n - k$ zeroes appended. Then $v^*A[k]v = w^*Aw$ is a positive real number, proving the claim.

Conversely, we now show that *if $A^* = A$ and $\det(A[k]) \in \mathbb{R}^+$ for $1 \leq k \leq n$, then A is positive definite*. We use induction on $n \geq 1$. For $n = 1$, it is immediate that $A \in M_1(\mathbb{C})$ is positive definite iff $A(1, 1) \in \mathbb{R}^+$ iff $\det(A[1]) = A(1, 1) \in \mathbb{R}^+$. Fix $n > 1$ and assume the result is known for matrices in $M_{n-1}(\mathbb{C})$. Fix a matrix $B \in M_n(\mathbb{C})$ such that $B^* = B$ and $\det(B[k]) \in \mathbb{R}^+$ for $1 \leq k \leq n$. Let $A = B[n-1] \in M_{n-1}(\mathbb{C})$ be obtained from B by deleting the last row and column. Note $A^* = A$ and $\det(A[k]) = \det(B[k]) \in \mathbb{R}^+$ for $1 \leq k \leq n-1$, so the induction hypothesis shows that A is positive definite. Let the eigenvalues of A be the positive real numbers $a_1 \leq a_2 \leq \cdots \leq a_{n-1}$, and let the eigenvalues of B be the real numbers $b_1 \leq b_2 \leq \cdots \leq b_n$. By the interlacing theorem in §7.14, we know

$b_1 \leq a_1 \leq b_2 \leq a_2 \leq \cdots \leq a_{n-1} \leq b_n$. Thus, every eigenvalue of B , except possibly b_1 , is $\geq a_1$ and hence is positive. On the other hand, $\det(B) = b_1 b_2 \cdots b_n > 0$ by hypothesis. Since $b_2, \dots, b_n > 0$, it follows that $b_1 > 0$ as well. Since B is normal and all its eigenvalues are positive real numbers, B is positive definite. This completes the induction.

7.16 Summary

We now summarize some facts about Hermitian, positive definite, unitary, and normal matrices.

1. *Definitions.* Given any complex matrix A , the *conjugate-transpose* A^* is the matrix obtained by transposing A and replacing each entry by its complex conjugate. A matrix $A \in M_n(\mathbb{C})$ is *Hermitian* iff $A^* = A$; A is *unitary* iff $A^* = A^{-1}$; A is *normal* iff $AA^* = A^*A$; A is *positive definite* iff v^*Av is a positive real number for all nonzero $v \in \mathbb{C}^n$; A is *positive semidefinite* iff v^*Av is a nonnegative real number for all $v \in \mathbb{C}^n$.
2. *Properties of Conjugate-Transpose.* Whenever the matrix operations are defined, the following identities hold: $(A+B)^* = A^* + B^*$; $(cA)^* = \bar{c}(A^*)$; $(AB)^* = B^*A^*$; $(A^*)^* = A$; $(A^{-1})^* = (A^*)^{-1}$; $\text{tr}(A^*) = \overline{\text{tr}(A)}$; $\det(A^*) = \overline{\det(A)}$; $\langle Av, w \rangle = w^*Av = \langle v, A^*w \rangle$.
3. *Hermitian Matrices.* The following conditions on a matrix $A \in M_n(\mathbb{C})$ are equivalent: A is Hermitian; $A = A^*$; v^*Av is real for all $v \in \mathbb{C}^n$; A is normal with all real eigenvalues. Real linear combinations of Hermitian matrices are Hermitian, as are products of *commuting* Hermitian matrices. Every matrix $B \in M_n(\mathbb{C})$ can be written uniquely in the form $B = X + iY$, where X and Y are Hermitian matrices; X and Y commute iff B is normal. This is called the *Cartesian or Hermitian decomposition* of B .
4. *Positive Definite Matrices.* The following conditions on a matrix $A \in M_n(\mathbb{C})$ are equivalent: A is positive definite; v^*Av is real and positive for all nonzero $v \in \mathbb{C}^n$; A is normal with all positive real eigenvalues; $A^* = A$ and for all $1 \leq k \leq n$, the matrix consisting of the first k rows and columns of A has positive determinant. Positive linear combinations of positive definite matrices are positive definite.
5. *Unitary Matrices.* The following conditions on a matrix $U \in M_n(\mathbb{C})$ are equivalent: U is unitary; $UU^* = I$; $U^*U = I$; $U^{-1} = U^*$; U^* is unitary; U^{-1} is unitary; U^T is unitary; the rows of U are orthonormal vectors in \mathbb{C}^n ; the columns of U are orthonormal vectors in \mathbb{C}^n ; $\langle Uv, Uw \rangle = \langle v, w \rangle$ for all $v, w \in \mathbb{C}^n$; U preserves inner products; $\|Uv\| = \|v\|$ for all $v \in \mathbb{C}^n$; U is length-preserving; U is the matrix of a length-preserving linear map relative to an orthonormal basis; U is normal with all eigenvalues having modulus 1 in \mathbb{C} . For U unitary, we have $|\det(U)| = 1$ in \mathbb{C} . Unitary similarity is an equivalence relation on square matrices that corresponds to an orthonormal change of basis on the underlying vector space.
6. *Unitary Triangularization and Diagonalization.* A matrix $A \in M_n(\mathbb{C})$ is *unitarily triangulable* iff $U^*AU = U^{-1}AU$ is upper-triangular for some unitary matrix U ; A is *unitarily diagonalizable* iff U^*AU is diagonal for some unitary matrix U . *Schur's theorem* says every matrix $A \in M_n(\mathbb{C})$ can be unitarily triangularized. The *spectral theorem* says A can be unitarily diagonalized iff A is normal iff there

exist n orthonormal eigenvectors for A . In this case, the eigenvectors are the columns of the diagonalizing matrix U , and the diagonal entries of U^*AU are the eigenvalues of A . The spectral theorem applies, in particular, to Hermitian, positive definite, and unitary matrices (which are all normal).

7. *Simultaneous Triangularization and Diagonalization.* A family $\{A_i : i \in I\}$ of $n \times n$ matrices can be *simultaneously* triangularized (resp. diagonalized) iff there exists a single invertible matrix S (independent of $i \in I$) such that $S^{-1}A_iS$ is upper-triangular (resp. diagonal) for all $i \in I$. A *commuting* family of $n \times n$ complex matrices can be simultaneously unitarily triangularized. A family of *normal* matrices can be simultaneously unitarily diagonalized iff the matrices in the family commute.
8. *k 'th Roots of Matrices.* For any even $k \geq 1$, every positive definite matrix A has a unique positive definite k 'th root, which is a matrix B such that $B^k = A$. Similarly, for any odd $k \geq 1$, every Hermitian matrix A has a unique Hermitian k 'th root. In each case, $B = p(A)$ for some polynomial p ; we can find B by unitarily diagonalizing A and taking positive (resp. real) k 'th roots of the positive (resp. real) diagonal entries.
9. *Polar Decomposition.* For every matrix $A \in M_n(\mathbb{C})$, there exist a positive semidefinite matrix R and a unitary matrix U such that $A = UR$. $R = \sqrt{A^*A}$ is always uniquely determined by A ; U is unique if A is invertible. U and R commute iff A is normal. There is a dual decomposition $A = R'U'$ even when A is not normal.
10. *Singular Value Decomposition.* Every matrix $A \in M_n(\mathbb{C})$ can be written in the form $A = PDQ^*$, where P and Q are unitary and D is a diagonal matrix with nonnegative real entries. The entries of D are uniquely determined by A (up to rearrangement); they are called the *singular values of A* . This result says that every linear transformation of \mathbb{C}^n is the composition of an isometry, then a nonnegative rescaling of orthonormal axes, then another isometry.
11. *Characterizations of Normality.* The following conditions on $A \in M_n(\mathbb{C})$ are equivalent: A is normal; $AA^* = A^*A$; A is unitarily diagonalizable; A has n orthonormal eigenvectors; the unique Hermitian matrices X and Y in the Cartesian decomposition $A = X + iY$ commute; $A^* = p(A)$ for some polynomial p ; the matrices U and R in the polar decomposition $A = UR$ commute.
12. *Interlacing Theorem.* Given $B \in M_n(\mathbb{C})$ with $B^* = B$, let A be B with the last row and column erased. If B has eigenvalues $b_1 \leq \dots \leq b_n$ and A has eigenvalues $a_1 \leq \dots \leq a_{n-1}$, then $b_1 \leq a_1 \leq b_2 \leq a_2 \leq \dots \leq b_{n-1} \leq a_{n-1} \leq b_n$.

7.17 Exercises

1. Prove the properties of complex conjugation stated in the introduction to this chapter.
2. Let $A = \begin{bmatrix} 1-i & 0 & 2+2i \\ 3 & -1+3i & 5i \end{bmatrix}$ and $B = \begin{bmatrix} i & -2+i \\ 1 & 1+3i \end{bmatrix}$. Compute A^* , B^* , AA^* , A^*A , $(BA)^*$, A^*B^* , B^{-1} , $(B^*)^{-1}$, and $(B^{-1})^*$.
3. For $A \in M_{m,n}(\mathbb{C})$, define $\bar{A} \in M_{m,n}(\mathbb{C})$ to be the matrix with i, j -entry $\overline{A(i,j)}$

for $i \in [m]$ and $j \in [n]$. Which of the eight properties of A^* in §7.1 have analogues for \overline{A} ? Prove your answers.

4. Let Q_A be the quadratic form associated with a matrix $A \in M_n(\mathbb{C})$ (see §7.1).
 (a) Prove: for all $A \in M_n(\mathbb{R})$, there exists a symmetric matrix $B \in M_n(\mathbb{R})$ with $Q_A = Q_B$. (b) Does part (a) hold if we replace \mathbb{R} by \mathbb{C} ? Explain. (c) Given $A \in M_n(\mathbb{C})$, must there exist a Hermitian $B \in M_n(\mathbb{C})$ with $Q_A = Q_B$?
5. Decide whether each matrix is Hermitian, unitary, positive definite, or normal (select all that apply, and explain): (a) I_n ; (b) $0 \in M_n(\mathbb{C})$;
 (c) $\begin{bmatrix} 0 & i \\ -i & 0 \end{bmatrix}$; (d) $\begin{bmatrix} 5 & -\sqrt{3} \\ -\sqrt{3} & 7 \end{bmatrix}$; (e) $\begin{bmatrix} 7 & 1 \\ 3 & 3 \end{bmatrix}$; (f) $\begin{bmatrix} 1/3 & 2/3 & 2/3 \\ 2/3 & -2/3 & 1/3 \\ 2/3 & 1/3 & -2/3 \end{bmatrix}$;
 (g) $\begin{bmatrix} -1+i & -4-2i & 6 \\ -4-2i & 3 & 2-2i \\ 6 & 2-2i & -2-i \end{bmatrix}$; (h) $\begin{bmatrix} 3 & i & 0 \\ i & 3 & i \\ 0 & i & 3 \end{bmatrix}$.
6. For a 2×2 matrix $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ with $a, b, c, d \in \mathbb{C}$, find algebraic conditions on a, b, c, d that are equivalent to A being: (a) Hermitian; (b) unitary; (c) positive definite; (d) positive semidefinite; (e) normal.
7. For a fixed $t \in \mathbb{C}$, call a matrix $A \in M_n(\mathbb{C})$ *t-Hermitian* iff $A^* = tA$. (a) Prove that the set of t -Hermitian matrices is a real subspace of $M_n(\mathbb{C})$. (b) Suppose A is t -Hermitian, B is s -Hermitian, and $AB = BA$. Prove AB is st -Hermitian.
8. (a) Prove: if $A \in M_n(\mathbb{C})$ is skew-Hermitian, then the associated quadratic form Q takes values in $\{ib : b \in \mathbb{R}\}$. (b) Is the converse of (a) true? Prove or give a counterexample. (c) Must a matrix with all pure imaginary eigenvalues be skew-Hermitian? Explain.
9. Find the Hermitian decomposition of each matrix: (a) $\begin{bmatrix} 2 & 5 \\ -1 & 3 \end{bmatrix}$;
 (b) $\begin{bmatrix} 1+3i & 2-5i \\ 3+4i & -1-2i \end{bmatrix}$; (c) $\begin{bmatrix} 3 & 7 \\ 7 & 3 \end{bmatrix}$; (d) $\begin{bmatrix} 2i & 1+3i \\ -1+3i & -5i \end{bmatrix}$;
 (e) $\begin{bmatrix} 0 & i & 2i \\ -i & 0 & 3i \\ 2i & 3i & 4 \end{bmatrix}$.
10. Given nonzero $v \in \mathbb{C}^n$, let $A = I - (2/v^*v)(vv^*) \in M_n(\mathbb{C})$. (a) Compute A for $v = (1, 2)$ and $v = (1, 7, 5, 5)$. (b) Show that A is Hermitian and unitary. (c) Show that $Av = -v$ and $Aw = w$ for all $w \in \mathbb{C}^n$ with $\langle w, v \rangle = 0$. (d) Describe all diagonal matrices unitarily similar to A .
11. Let F be a field in which $1_F + 1_F \neq 0_F$. (a) Prove: for all $A \in M_n(F)$, there exist unique $B, C \in M_n(F)$ with $A = B + C$, $B^T = B$, and $C^T = -C$. (b) Does (a) hold when $1_F + 1_F = 0_F$? Explain.
12. Let $A \in M_n(\mathbb{C})$ be positive definite. (a) Show A^T is positive definite. (b) Show A^{-1} is positive definite. (c) For $k \in \mathbb{N}^+$, must A^k be positive definite? Explain.
13. Give an example of a matrix A of size $n \geq 2$ satisfying the indicated properties, or explain why no such example exists. (a) A is positive definite, but $A(i, j)$ is a negative real number for all $i, j \in [n]$; (b) A is positive definite, but $A(i, j)$ is pure imaginary for all $i \neq j$ in $[n]$; (c) A is positive semidefinite and negative semidefinite and nonzero; (d) A is negative definite, but $\det(A) > 0$; (e) all eigenvalues of A are real and negative, but A is not negative definite.

14. Suppose $A \in M_n(\mathbb{C})$ satisfies $v^*Av \in \mathbb{R}^+$ for all v in a fixed basis of \mathbb{C}^n . Give an example to show that A need not be positive definite.
15. Suppose $A \in M_n(\mathbb{R})$ satisfies $A^T = A$ and $v^TAv \geq 0$ for all nonzero $v \in \mathbb{R}^n$.
 (a) Prove A is positive semidefinite. (b) Give a specific example demonstrating that the conclusion of (a) fails without the hypothesis $A^T = A$.
16. Use the determinant criterion in §7.15 to find all $c \in \mathbb{C}$ such that the following matrices are positive definite. (a) $A = \begin{bmatrix} 2 & c \\ c & 3 \end{bmatrix}$; (b) $A = \begin{bmatrix} 4 & c & 0 \\ c & 4 & c \\ 0 & c & 4 \end{bmatrix}$;
 (c) $A = \begin{bmatrix} c & 1 & 2 & 0 \\ 1 & c & 1 & 2 \\ 2 & 1 & c & 1 \\ 0 & 2 & 1 & c \end{bmatrix}$.
17. (a) Suppose $A \in M_n(\mathbb{C})$ satisfies $A^* = A$ and $\det(A[k]) \geq 0$ for $1 \leq k \leq n$. Give an example to show that A need not be positive semidefinite. (b) Now assume $A^* = A$, $\det(A[k]) > 0$ for $1 \leq k < n$, and $\det(A) \geq 0$. Prove A is positive semidefinite.
18. State and prove a determinant characterization of negative definite matrices.
19. Suppose $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ has continuous second partial derivatives, and f has a local minimum at $(a, b) \in \mathbb{R}^2$. (a) Prove that $A = \begin{bmatrix} f_{xx}(a, b) & f_{xy}(a, b) \\ f_{yx}(a, b) & f_{yy}(a, b) \end{bmatrix}$ is positive semidefinite. Deduce that $f_{xx}f_{yy} - f_{xy}^2 \geq 0$ and $f_{xx}, f_{yy} \geq 0$ at (a, b) . [Hint: For any nonzero $v = (v_1, v_2) \in \mathbb{R}^2$, study $g(t) = f((a, b) + t(v_1, v_2))$ for $t \in \mathbb{R}$.] (b) State and prove a version of (a) for a local maximum of f . (c) Generalize (a) and (b) to functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$.
20. How many diagonal unitary matrices are there in $M_n(\mathbb{R})$?
21. Give an example of a unitary matrix U with each property, or explain why this cannot be done. (a) $U^T = -U$; (b) $U(1, 1) = 2$; (c) $U^* = U \neq I$; (d) U is upper-triangular and not diagonal; (e) U is symmetric with all entries nonzero.
22. Let U be any 2×2 real unitary matrix. (a) Show that the first row of U has the form $(\cos \theta, \sin \theta)$ for some $\theta \in [0, 2\pi)$. (b) Given the first row of U as in (a), find all possible second rows of U .
23. Let Z be any subspace of \mathbb{C}^n . Prove there exists an orthonormal basis of Z by induction on $\dim(Z)$. (Use the claim in §7.7.)
24. Let W be a subspace of \mathbb{C}^n with orthonormal basis (w_1, \dots, w_k) . Define the *orthogonal complement* W^\perp to be the set of $z \in \mathbb{C}^n$ such that $\langle z, w \rangle = 0$ for all $w \in W$. (a) Prove $z \in W^\perp$ iff $\langle z, w_i \rangle = 0$ for $1 \leq i \leq k$. (b) Prove W^\perp is a subspace of \mathbb{C}^n . (c) Prove $W \cap W^\perp = \{0\}$. (d) Prove $\dim(W) + \dim(W^\perp) = n$. (e) Conclude that $\mathbb{C}^n = W \oplus W^\perp$.
25. Let W be a subspace of \mathbb{C}^n with orthonormal basis (w_1, \dots, w_k) . We can extend this list to a basis $(w_1, \dots, w_k, x_{k+1}, \dots, x_n)$ of \mathbb{C}^n . For $k+1 \leq j \leq n$, define

$$y_j = x_j - \sum_{r=1}^{j-1} (w_r^* x_j) w_r, \quad w_j = y_j / \|y_j\|.$$

Use induction to prove these facts: (a) each y_j is nonzero; (b) $\langle y_j, w_s \rangle = 0$ for all $s < j$; (c) (w_1, \dots, w_j) spans the same subspace as $(w_1, \dots, w_k, x_{k+1}, \dots, x_j)$

- for $k + 1 \leq j \leq n$; (d) (w_1, \dots, w_n) is an orthonormal basis for \mathbb{C}^n . So, *every orthonormal list in \mathbb{C}^n can be extended to an orthonormal basis of \mathbb{C}^n* .
26. Fix $A \in M_n(\mathbb{C})$. (a) Prove: if there exists $U \in U_n(\mathbb{C})$ such that A is similar to U , then A^{-1} is similar to A^* . (b) Prove or disprove the converse of (a).
 27. Find all $A \in M_n(\mathbb{C})$ such that A is the only matrix unitarily similar to A .
 28. Given $A \in M_n(\mathbb{C})$ and $f \in S_n$, define $B \in M_n(\mathbb{C})$ by $B(i, j) = A(f(i), f(j))$ for all $i, j \in [n]$. Prove A and B are unitarily similar.
 29. Given any subgroup H of $\mathrm{GL}_n(\mathbb{C})$, call two matrices $A, B \in M_n(\mathbb{C})$ *H -similar* iff $B = S^{-1}AS$ for some $S \in H$. (a) Show that H -similarity is an equivalence relation on $M_n(\mathbb{C})$. (b) Give an example of subgroups $H \neq K$ in $\mathrm{GL}_n(\mathbb{C})$ such that H -similarity coincides with K -similarity.
 30. With as little calculation as possible, explain why each pair of matrices cannot be unitarily similar. (a) $A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$, $B = \begin{bmatrix} -1 & -8 \\ 1 & 5 \end{bmatrix}$; (b) $A = \begin{bmatrix} 0.6 & -0.8 \\ 0.8 & 0.6 \end{bmatrix}$, $B = \begin{bmatrix} -3.4 & -1.6 \\ 10.4 & 4.6 \end{bmatrix}$; (c) $A = \begin{bmatrix} 2 & 2 \\ 2 & 5 \end{bmatrix}$, $B = \begin{bmatrix} 9 & 2 \\ -12 & -2 \end{bmatrix}$.
 31. Suppose $e_1 \cdots e_k$ is any finite sequence of 1's and *'s. Prove: for all $A, B \in M_n(\mathbb{C})$, if A is unitarily similar to B , then $\mathrm{tr}(A^{e_1} A^{e_2} \cdots A^{e_k}) = \mathrm{tr}(B^{e_1} B^{e_2} \cdots B^{e_k})$. (For example, if $e_1 e_2 e_3 e_4 = *1**$, the condition states that $\mathrm{tr}(A^* A A^* A^*) = \mathrm{tr}(B^* B B^* B^*)$.) Conversely, it can be shown that if the equality of traces holds for all sequences $e_1 \cdots e_k$, then A and B are unitarily similar.
 32. (a) Prove that any normal matrix $A \in M_n(\mathbb{C})$ is unitarily similar to A^T . (b) Is every $A \in M_n(\mathbb{C})$ unitarily similar to A^T ?
 33. For each matrix A , find a unitary U and an upper-triangular T with $T = U^*AU$.
 - (a) $A = \begin{bmatrix} 2.5 & -0.5 \\ 0.5 & 3.5 \end{bmatrix}$; (b) $A = \begin{bmatrix} 5 & -108 & -36 \\ 39 & -127 & -26 \\ -36 & 72 & -25 \end{bmatrix}$;
 - (c) $A = \begin{bmatrix} 4 & 1 & 0 & 0 \\ 1 & 4 & 1 & 0 \\ 0 & 1 & 4 & 1 \\ 0 & 0 & 1 & 4 \end{bmatrix}$.
 34. Define $A \in M_5(\mathbb{C})$ by $A(i, j) = i$ for $i, j \in [5]$. (a) Find a matrix $S \in \mathrm{GL}_5(\mathbb{C})$ such that $S^{-1}AS$ is diagonal. (b) Find a matrix $U \in U_5(\mathbb{C})$ such that U^*AU is upper-triangular. (c) Is there a matrix $V \in U_n(\mathbb{C})$ such that V^*AV is diagonal? Why?
 35. Give an example of a matrix $A \in M_2(\mathbb{R})$ for which there does not exist any unitary $U \in M_2(\mathbb{R})$ with U^*AU upper-triangular.
 36. (a) Suppose $A, B \in M_n(\mathbb{C})$ are unitarily similar. Prove $\sum_{i,j \in [n]} |A(i, j)|^2 = \sum_{i,j \in [n]} |B(i, j)|^2$. [Hint: What is the trace of AA^* ?] (b) Suppose A is unitarily similar to two upper-triangular matrices T_1 and T_2 . Prove $\sum_{i < j} |T_1(i, j)|^2 = \sum_{i < j} |T_2(i, j)|^2$. (c) Give an example of upper-triangular matrices T_1 and $T_2 \in M_3(\mathbb{C})$ such that T_1 and T_2 are unitarily similar, $T_1(i, i) = T_2(i, i)$ for $i = 1, 2, 3$, but $T_1 \neq T_2$. (d) Give an example of upper-triangular matrices T_1 and T_2 with the same diagonal such that $\sum_{i < j} |T_1(i, j)|^2 = \sum_{i < j} |T_2(i, j)|^2$, but T_1 and T_2 are not unitarily similar.
 37. Let $A \in M_n(\mathbb{C})$ have eigenvalues $c_1, \dots, c_n \in \mathbb{C}$. Prove A is normal iff $\sum_{i,j \in [n]} |A(i, j)|^2 = \sum_{i=1}^n |c_i|^2$.

38. Which pairs of matrices can be simultaneously unitarily triangularized? Which can be simultaneously unitarily diagonalized? Decide without calculating any eigenvalues. (a) $A = \begin{bmatrix} 1 & 4 \\ 3 & 1 \end{bmatrix}$, $B = \begin{bmatrix} 0 & 12 \\ 9 & 0 \end{bmatrix}$; (b) $A = \begin{bmatrix} 3 & -1 \\ 0 & 2 \end{bmatrix}$, $B = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix}$; (c) $A = \begin{bmatrix} a & b \\ b & a \end{bmatrix}$, $B = \begin{bmatrix} c & d \\ d & c \end{bmatrix}$ where $a, b, c, d \in \mathbb{R}$.
39. (a) For even $k \in \mathbb{N}^+$, show that for each positive definite (resp. semidefinite) $A \in M_n(\mathbb{C})$, there exists a unique positive definite (resp. semidefinite) $B \in M_n(\mathbb{C})$ with $B^k = A$. (b) For odd $k \in \mathbb{N}^+$, show that for each Hermitian $A \in M_n(\mathbb{C})$, there exists a unique Hermitian $B \in M_n(\mathbb{C})$ with $B^k = A$.
40. Compute a numerical approximation to the unique positive definite square root of each positive definite matrix. (a) $\begin{bmatrix} 5 & -1 \\ -1 & 5 \end{bmatrix}$; (b) $\begin{bmatrix} 3 & 2 & 1 \\ 2 & 3 & 2 \\ 1 & 2 & 3 \end{bmatrix}$; (c) $\begin{bmatrix} 2 & 1+i \\ 1-i & 2 \end{bmatrix}$.
41. Prove $A \in M_n(\mathbb{C})$ is positive definite iff there is an invertible $B \in M_n(\mathbb{C})$ with $A = B^*B$.
42. (a) Give an example of a matrix $A \in M_2(\mathbb{C})$ such that $B^2 = A$ has no solution $B \in M_2(\mathbb{C})$. (b) For each $n > 2$, find $A \in M_n(\mathbb{C})$ with no square root in $M_n(\mathbb{C})$.
43. Fix $n > 1$. (a) Show that $I_n \in M_n(\mathbb{C})$ has infinitely many square roots in $M_n(\mathbb{C})$. (b) Show that $0 \in M_n(\mathbb{C})$ has infinitely many square roots in $M_n(\mathbb{C})$.
44. (a) Let $T \in M_2(\mathbb{C})$ be upper-triangular. Find all upper-triangular $U \in M_2(\mathbb{C})$ with $U^2 = T$. (b) Repeat (a) for $M_3(\mathbb{C})$. (c) Do there exist matrices T, U such that $U^2 = T$, T is upper-triangular, but U is not upper-triangular?
45. Let $T \in M_n(\mathbb{C})$ be upper-triangular with distinct nonzero diagonal entries. Prove there exist exactly 2^n upper-triangular $U \in M_n(\mathbb{C})$ with $U^2 = T$.
46. Prove that for $n > 1$, a matrix $A \in M_n(\mathbb{C})$ cannot have a finite odd number of square roots in $M_n(\mathbb{C})$.
47. Find (with proof) an example of a matrix $A \in M_2(\mathbb{C})$ that has exactly two square roots in $M_2(\mathbb{C})$.
48. (a) Suppose $A \in M_n(\mathbb{C})$ has eigenvalues c_1, \dots, c_n , $B \in M_n(\mathbb{C})$ has eigenvalues d_1, \dots, d_n , and $AB = BA$. Prove there exists $f \in S_n$ such that $A + B$ has eigenvalues $c_1 + d_{f(1)}, \dots, c_n + d_{f(n)}$. (b) With the setup in (a), what can you say about the eigenvalues of AB ? (c) Give an example to show that (a) can fail without the hypothesis $AB = BA$.
49. For each matrix A , find a polynomial $p \in \mathbb{C}[x]$ with $A^* = p(A)$, or explain why this is impossible. (a) $A = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 3+4i & 0 \\ 0 & 0 & 2 \end{bmatrix}$; (b) $A = \begin{bmatrix} 5 & 5 \\ 2 & 2 \end{bmatrix}$; (c) $A = \begin{bmatrix} 5 & 2 \\ -2 & 5 \end{bmatrix}$; (d) $A = \begin{bmatrix} 1 & 1+2i & 1 & -1+2i \\ -1-2i & 1 & -1+2i & -1 \\ 1 & 1-2i & 1 & -1-2i \\ 1-2i & -1 & 1+2i & 1 \end{bmatrix}$.
50. Suppose $A \in M_n(\mathbb{C})$ satisfies $A^{-1} = p(A)$ for some $p \in \mathbb{C}[x]$. Must A be normal?
51. (a) Suppose $A \in M_n(\mathbb{C})$ is normal. Must $\overline{A} = p(A)$ for some $p \in \mathbb{C}[x]$? (b) Suppose $\overline{A} = p(A)$ for some $p \in \mathbb{C}[x]$. Must A be normal?

52. (a) Suppose $A \in M_n(\mathbb{C})$ is normal. Must $A^T = p(A)$ for some $p \in \mathbb{C}[x]$?
 (b) Suppose $A^T = p(A)$ for some $p \in \mathbb{C}[x]$. Must A be normal?
53. Find a polar decomposition $A = UR$ for each matrix below: (a) $A = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$;
 (b) $A = \begin{bmatrix} 2 & 2 \\ 2 & 2 \end{bmatrix}$; (c) $A = \begin{bmatrix} -5 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -7 \end{bmatrix}$; (d) $\begin{bmatrix} 0 & 0 & 0 \\ 2 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$.
54. Find a singular value decomposition and the singular values for each matrix in Exercise 53.
55. Prove the following “abstract” version of the singular value decomposition: given a linear map $T : V \rightarrow W$ between two n -dimensional complex inner product spaces, there exist an orthonormal basis $X = (x_1, \dots, x_n)$ for V and an orthonormal basis $Y = (y_1, \dots, y_n)$ for W and nonnegative real numbers c_1, \dots, c_n such that $T(x_i) = c_i y_i$ for $1 \leq i \leq n$.
56. Prove the singular value decomposition for rectangular matrices: given $A \in M_{m,n}(\mathbb{C})$, there exist unitary matrices $P \in M_m(\mathbb{C})$ and $Q \in M_n(\mathbb{C})$ and a matrix $D \in M_{m,n}(\mathbb{C})$ with $D(i,i) \geq 0$ for all i and $D(i,j) = 0$ for $i \neq j$, such that $A = PDQ^*$.
57. Verify (7.4).
58. In step 2 of §7.14, assume $a_{j-1} < a_j = \dots = a_{j+s} < a_{j+s+1}$. (a) Compute the sign of $\lim_{\epsilon \rightarrow 0^+} \chi_B(a_j - \epsilon)/\epsilon^s$. (b) Compute the sign of $\lim_{\epsilon \rightarrow 0^+} \chi_B(a_{j+s} + \epsilon)/\epsilon^s$. (c) Use (a) and (b) to complete the proof of step 2.
59. For the matrix $A = \begin{bmatrix} 1 & 4 & 0 & 0 \\ 4 & 7 & -1 & -5 \\ 0 & -1 & 1 & 0 \\ 0 & -5 & 0 & 5 \end{bmatrix}$, verify the interlacing theorem in §7.14 by computing the eigenvalues of $A[k]$ for $k = 1, 2, 3, 4$.
60. True or false? Explain each answer. (Assume matrices are in $M_n(\mathbb{C})$ unless otherwise stated.) (a) Every diagonal matrix is normal. (b) A is positive semidefinite iff $-A$ is negative semidefinite. (c) For all $A \in M_{m,n}(\mathbb{C})$, AA^* is positive definite. (d) For all $A \in M_{m,n}(\mathbb{C})$, A^*A is positive semidefinite. (e) Every real symmetric matrix is unitarily diagonalizable. (f) The product of two Hermitian matrices is always Hermitian. (g) The product of two positive definite matrices is always positive definite. (h) The product of two unitary matrices is always unitary. (i) The square of a Hermitian matrix is always positive semidefinite. (j) The inverse of a negative definite matrix exists and must be negative definite. (k) The inverse of an invertible normal matrix must be normal. (l) If A is normal, then $p(A)$ is normal for all $p \in \mathbb{C}[x]$. (m) Every positive definite matrix has a unique square root. (n) For fixed n , I_n is the only matrix that is both unitary and Hermitian. (o) A normal matrix with all entries in \mathbb{R}^+ must be positive definite. (p) If $A \in M_2(\mathbb{C})$ has $A(1,1) \geq 0$ and $\det(A) \geq 0$, then A must be positive semidefinite. (q) If every complex eigenvalue of a matrix has modulus 1, the matrix must be unitary.

Jordan Canonical Forms

Let V be an n -dimensional vector space over a field F . Given any linear operator $T : V \rightarrow V$ and any ordered basis X of V , recall from Chapter 6 that there is a matrix $A = [T]_X$ in $M_n(F)$ that represents T relative to the basis X . If we switch from X to another ordered basis Y , the matrix A is replaced by a similar matrix of the form $[T]_Y = P^{-1}AP$, where $P \in M_n(F)$ is some invertible matrix (specifically, P is the matrix of id_V relative to the input basis Y and output basis X).

A fundamental question in linear algebra is how to pick the ordered basis Y to make the matrix $[T]_Y$ as simple as possible. From the viewpoint of matrix computations, the simplest $n \times n$ matrices are *diagonal* matrices D , which satisfy $D(i, j) = 0_F$ for all $i \neq j$. A linear map T is called *diagonalizable* iff there exists an ordered basis Y such that $[T]_Y$ is a diagonal matrix. Regrettably, for $n > 1$, not all linear maps on V can be diagonalized. For example, suppose T is a *nilpotent* linear operator, which means that $T^k = T \circ T \circ \dots \circ T$ (k factors) $= 0$ for some positive integer k . Suppose $[T]_Y$ were a diagonal matrix D . On one hand, $[T^k]_Y = [0]_Y = 0$. On the other hand, $[T^k]_Y = [T]_Y^k = D^k$. Since D is diagonal, $D^k = 0$ forces every diagonal entry of D to be zero, so that $D = 0$ and $T = 0$. But for $n > 1$, one can find many examples of nonzero nilpotent linear maps (as we will see below). These maps cannot be diagonalized.

In this chapter, we will prove that linear maps on a *complex* vector space can always be represented by certain “nearly-diagonal” matrices called *Jordan canonical forms*. More precisely, for any field F , any $c \in F$, and any $m \in \mathbb{N}^+$, define the *Jordan block*

$$J(c; m) = \begin{bmatrix} c & 1 & 0 & \cdots & 0 \\ 0 & c & 1 & \cdots & 0 \\ 0 & 0 & c & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 1 \\ 0 & 0 & 0 & \cdots & c \end{bmatrix}_{m \times m}.$$

This matrix has m copies of c on its diagonal, $m - 1$ ones on the next higher diagonal, and zeroes elsewhere. When $m = 1$, $J(c; 1)$ is the 1×1 matrix $[c]$. Next, define a *Jordan canonical form* to be any matrix \mathbf{J} that has the block-diagonal structure

$$\mathbf{J} = \begin{bmatrix} J(c_1; m_1) & 0 & \cdots & 0 \\ 0 & J(c_2; m_2) & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & J(c_s; m_s) \end{bmatrix} \quad (8.1)$$

for some $s, m_1, \dots, m_s \in \mathbb{N}^+$ and $c_1, \dots, c_s \in F$ (the m_j 's and c_j 's need not all be distinct). Writing $\text{blk-diag}(A_1, \dots, A_s)$ to denote a block-diagonal matrix with diagonal blocks A_1, \dots, A_s , we have $\mathbf{J} = \text{blk-diag}(J(c_1; m_1), \dots, J(c_s; m_s))$.

The **Jordan canonical form theorem** states that *given any linear map T on a finite-dimensional complex vector space V , there is an ordered basis X of V such that the matrix $\mathbf{J} = [T]_X$ is a Jordan canonical form. Furthermore, if Y is any ordered basis such that*

$\mathbf{J}' = [T]_Y$ is also a Jordan canonical form, then \mathbf{J}' is obtained from \mathbf{J} by rearranging the Jordan blocks in (8.1). The theorem can also be phrased as a statement about matrices: given any $A \in M_n(\mathbb{C})$, there exists an invertible $P \in M_n(\mathbb{C})$ such that $\mathbf{J} = P^{-1}AP$ is a Jordan canonical form. If $\mathbf{J}' = Q^{-1}AQ$ is any Jordan canonical form similar to A , then \mathbf{J}' is obtained by rearranging the Jordan blocks of \mathbf{J} .

To prove the Jordan canonical form theorem, we first analyze the structure of nilpotent maps. Given an n -dimensional vector space V over any field F (not just the field \mathbb{C}) and any nilpotent linear map $T : V \rightarrow V$, we will prove that there exist unique integers $m_1 \geq m_2 \geq \dots \geq m_s > 0$ such that for some ordered basis X of V ,

$$[T]_X = \text{blk-diag}(J(0_F; m_1), J(0_F; m_2), \dots, J(0_F; m_s)).$$

This proof will be facilitated by the visual analysis of *partition diagrams* that represent the sequence (m_1, \dots, m_s) . Next, we discuss Fitting's lemma, which (roughly speaking) takes an arbitrary linear map $T : V \rightarrow V$ and breaks V into two pieces, such that T restricted to one piece is nilpotent, and T restricted to the other piece is an isomorphism. Combining this result with the classification of nilpotent maps, we will arrive at the Jordan canonical form theorem stated above.

After proving this theorem, we discuss methods for actually computing the Jordan canonical form of a specific linear map or matrix. Then we describe an application of this theory to the solution of systems of linear ordinary differential equations. Finally, we study a more abstract version of the Jordan canonical form theorem that is needed (for example) in the study of Lie algebras. We will prove that every linear map on a finite-dimensional complex vector space V can be written uniquely as the sum of a diagonalizable linear map and a nilpotent linear map that commute with each other.

8.1 Examples of Nilpotent Maps

We begin our analysis of nilpotent linear maps by constructing some specific examples of nonzero nilpotent maps. Let F be any field, and let V be an F -vector space with ordered basis $X = (x_1, \dots, x_n)$. Recall that we can define an F -linear map $T : V \rightarrow V$ by choosing any elements $y_1, \dots, y_n \in V$, declaring that $T(x_j) = y_j$ for $1 \leq j \leq n$, and then setting $T(\sum_{j=1}^n c_j x_j) = \sum_{j=1}^n c_j y_j$ for all $c_j \in F$. In other words, we can build a unique linear map by sending basis elements anywhere and then *extending by linearity*. Furthermore, writing $T(x_j) = y_j = \sum_{i=1}^n a_{ij} x_i$ for $a_{ij} \in F$, we know that the matrix $A = [T]_X$ has i, j -entry a_{ij} .

To illustrate this procedure, suppose $n = 5$ and we decide that $T(x_1) = 0$, $T(x_2) = x_1$, $T(x_3) = x_2$, $T(x_4) = x_3$, and $T(x_5) = x_4$. The following “arrow diagram” illustrates the action of T on the basis X :

$$0 \xleftarrow{T} x_1 \xleftarrow{T} x_2 \xleftarrow{T} x_3 \xleftarrow{T} x_4 \xleftarrow{T} x_5.$$

For a general vector $v = c_1 x_1 + c_2 x_2 + c_3 x_3 + c_4 x_4 + c_5 x_5$ with $c_j \in F$,

$$T(v) = c_1 \cdot 0 + c_2 x_1 + c_3 x_2 + c_4 x_3 + c_5 x_4.$$

Applying T again, we find that

$$\begin{aligned} T^2(v) &= T(T(v)) = 0 + c_3 x_1 + c_4 x_2 + c_5 x_3, \\ T^3(v) &= c_4 x_1 + c_5 x_2, \\ T^4(v) &= c_5 x_1, \\ T^5(v) &= 0. \end{aligned}$$

Since $T^5(v) = 0$ for all $v \in V$, $T^5 = 0$ and T is nilpotent. On the other hand, the powers T^k for $1 \leq k < 5$ are nonzero maps. In general, we say a nilpotent map T has *index of nilpotence* m iff m is the least positive integer with $T^m = 0$. The matrix of T relative to the ordered basis X is

$$\begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} = J(0; 5).$$

Generalizing this example, for any ordered basis $X = (x_1, \dots, x_n)$, we could define a linear map T by setting $T(x_1) = 0$ and $T(x_j) = x_{j-1}$ for $2 \leq j \leq n$:

$$0 \xleftarrow{T} x_1 \xleftarrow{T} x_2 \xleftarrow{T} x_3 \xleftarrow{T} \cdots \xleftarrow{T} x_{n-1} \xleftarrow{T} x_n.$$

Since $T(x_1) = 0$ and $T(x_j) = 1x_{j-1} + \sum_{k \neq j-1} 0x_k$ for $j > 1$, we see that $[T]_X$ has ones in the $(j-1, j)$ -positions (for $2 \leq j \leq n$) and zeroes elsewhere. In other words, $[T]_X = J(0; n)$. Moreover, $T(\sum_{k=1}^n c_k x_k) = \sum_{k=2}^n c_k x_{k-1}$ for all $c_k \in F$. Iterating this formula gives $T^i(\sum_{k=1}^n c_k x_k) = \sum_{k=i+1}^n c_k x_{k-i}$ for all $i \geq 0$. So $T^i = 0$ for all $i \geq n$, and T is nilpotent of index n .

Next, consider an example where $n = 9$ and T acts on basis vectors as shown here:

$$\begin{aligned} 0 &\xleftarrow{T} x_1 \xleftarrow{T} x_2 \xleftarrow{T} x_3 \\ 0 &\xleftarrow{T} x_4 \xleftarrow{T} x_5 \\ 0 &\xleftarrow{T} x_6 \xleftarrow{T} x_7 \\ 0 &\xleftarrow{T} x_8 \\ 0 &\xleftarrow{T} x_9 \end{aligned} \tag{8.2}$$

One may check that $[T]_X = \text{blk-diag}(J(0; 3), J(0; 2), J(0; 2), J(0; 1), J(0; 1))$, $T \neq 0$, $T^2 \neq 0$, but $T^3 = 0$. So T is a nilpotent map of index 3.

For our next example, let $n = 4$ and define T on the basis by $T(x_1) = 0$, $T(x_2) = x_1$, $T(x_3) = 2x_1$, and $T(x_4) = 3x_1$. One sees that T is nilpotent of index 2, and

$$[T]_X = \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

This matrix is not a Jordan canonical form. However, if we let $Z = (x_1, x_2, x_3 - 2x_2, x_4 - 3x_2)$, then one checks that Z is an ordered basis of V and $[T]_Z = \text{blk-diag}(J(0; 2), J(0; 1), J(0; 1))$.

8.2 Partition Diagrams

Intuitively, we want to show that every nilpotent linear map on V can be described by an arrow diagram like (8.2), if we choose the right ordered basis for V . To discuss these arrow diagrams more precisely, we introduce the idea of *partition diagrams*.

A *partition* of an integer $n \in \mathbb{N}$ is a sequence $\mu = (\mu_1, \mu_2, \dots, \mu_s)$ where each μ_i is a positive integer, $\mu_1 \geq \mu_2 \geq \dots \geq \mu_s$, and $\mu_1 + \mu_2 + \dots + \mu_s = n$. We let $\ell(\mu)$ be the length

of the sequence μ . For example, $\mu = (3, 2, 2, 1, 1)$ is a partition of 9 with $\ell(\mu) = 5$. For any $c \in F$, define

$$J(c; \mu) = \text{blk-diag}(J(c; \mu_1), J(c; \mu_2), \dots, J(c; \mu_s)).$$

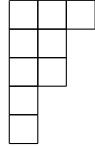
Given any partition μ , the *diagram* of μ is the set

$$D(\mu) = \{(i, j) \in \mathbb{N}^+ \times \mathbb{N}^+ : 1 \leq i \leq \ell(\mu), 1 \leq j \leq \mu_i\}.$$

We often visualize $D(\mu)$ by drawing a picture with a box in row i and column j for each $(i, j) \in D(\mu)$. For example, the diagram of $(3, 2, 2, 1, 1)$ is the set

$$\{(1, 1), (1, 2), (1, 3), (2, 1), (2, 2), (3, 1), (3, 2), (4, 1), (5, 1)\},$$

which is drawn as follows:



For all $k \geq 1$, let μ'_k be the number of boxes in the k 'th column of the diagram of μ . In our example, $\mu'_1 = 5$, $\mu'_2 = 3$, $\mu'_3 = 1$, and $\mu'_k = 0$ for all $k > 3$. In general, $\mu'_1 = \ell(\mu)$ and $\mu'_k = 0$ for all $k > \mu_1$.

Evidently, one can recover a partition μ from its diagram $D(\mu)$. In turn, the diagram $D(\mu)$ can be reconstructed from the sequence of column lengths (μ'_1, μ'_2, \dots) . Similarly, we can find the column lengths if we are given the sequence of partial sums $(\mu'_1, \mu'_1 + \mu'_2, \mu'_1 + \mu'_2 + \mu'_3, \dots)$, where the k 'th entry in the sequence counts the number of boxes in the first k columns. So, for any partitions μ and ν ,

$$\mu = \nu \text{ iff for all } k \geq 1, \mu'_1 + \dots + \mu'_k = \nu'_1 + \dots + \nu'_k. \quad (8.3)$$

This observation will be the key to establishing the uniqueness properties of Jordan canonical forms.

8.3 Partition Diagrams and Nilpotent Maps

Let V be an n -dimensional vector space over any field F . For each partition μ of n and each ordered basis $X = (x_1, \dots, x_n)$ of V , we can build a nilpotent linear map $T = T_{X, \mu}$ as follows. Fill the boxes in the diagram of μ with the elements of X , row by row, working from left to right in each row. Define the map T by sending each basis vector in the leftmost column to zero, sending every other basis vector to the element to its immediate left, and extending by linearity. For example, given $\mu = (5, 5, 3, 1, 1, 1)$ and $X = (x_1, \dots, x_{16})$, we first fill $D(\mu)$ as shown here:

x_1	x_2	x_3	x_4	x_5
x_6	x_7	x_8	x_9	x_{10}
x_{11}	x_{12}		x_{13}	
x_{14}				
x_{15}				
x_{16}				

Then $T_{X,\mu}$ is the unique linear map on V that sends basis vectors $x_1, x_6, x_{11}, x_{14}, x_{15}$, and x_{16} to 0, x_2 to x_1 , x_7 to x_6 , x_9 to x_8 , and so on. One readily checks that T is nilpotent of index 5, and $[T]_X = J(0; (5, 5, 3, 1, 1, 1))$.

Let us describe the preceding construction a bit more formally. Given a partition μ of n and an ordered basis $X = (x_1, \dots, x_n)$, we can view the *ordered* basis X as an *indexed set* $X = \{x(i, j) : (i, j) \in D(\mu)\}$ by setting $x(i, j) = x_{j+\mu_1+\mu_2+\dots+\mu_{i-1}}$ for all $(i, j) \in D(\mu)$. Intuitively, $x(i, j)$ is the basis vector placed in cell (i, j) of the diagram. We define $T = T_{X,\mu}$ on the basis X by setting $T(x(i, 1)) = 0_V$ for $1 \leq i \leq \ell(\mu)$ and $T(x(i, j)) = x(i, j - 1)$ for all $(i, j) \in D(\mu)$ with $j > 1$. Extending by linearity, we see that

$$T \left(\sum_{(i,j) \in D(\mu)} c(i, j)x(i, j) \right) = \sum_{i=1}^{\ell(\mu)} \sum_{j=2}^{\mu_i} c(i, j)x(i, j - 1) \quad (c(i, j) \in F). \quad (8.4)$$

Returning to the ordered basis (x_1, \dots, x_n) , we see that T sends $x_1, x_{\mu_1+1}, x_{\mu_1+\mu_2+1}$, etc., to zero, and T sends every other x_k to x_{k-1} . This observation is equivalent to the statement that $[T]_X = J(0; \mu)$.

Our goal is to prove that for every nilpotent linear map T on V , there exists a unique partition μ of $n = \dim(V)$ such that $[T]_X = J(0; \mu)$ for some ordered basis X of V . By the preceding remarks, this is equivalent to proving that for every nilpotent linear map T on V , there exists a unique partition μ and a (not necessarily unique) ordered basis X with $T = T_{X,\mu}$. Before giving this proof, we show how the partition diagram $D(\mu)$ encodes information about the image and null space of $T_{X,\mu}$ and its powers.

8.4 Computing Images via Partition Diagrams

Given any linear map $T : V \rightarrow V$, recall that the *image* of T (also called the *range* of T) is the subspace $\text{img}(T) = T[V] = \{T(v) : v \in V\}$. The *null space* of T (also called the *kernel* of T) is the subspace $\text{Null}(T) = \ker(T) = \{v \in V : T(v) = 0_V\}$. Given a nilpotent map of the special form $T = T_{X,\mu}$, let us see how to use the diagram of μ to compute the dimensions of $\text{img}(T^k)$ and $\text{Null}(T^k)$ for all $k \geq 1$. We will use $\mu = (5, 5, 3, 1, 1, 1)$ as a running example to illustrate the general discussion.

We begin by giving a visual representation of the formula (8.4) defining $T_{X,\mu}$. Represent a vector $v = \sum_{(i,j) \in D(\mu)} c(i, j)x(i, j)$ in V by filling each cell $(i, j) \in D(\mu)$ with the scalar $c(i, j) = c_{i,j} \in F$ that multiplies the basis vector $x(i, j)$ indexed by that cell. Applying T shifts all these coefficients one cell to the left, with coefficients $c(i, 1)$ falling off the left edge and new zero coefficients coming into the right cell of each row. For example,

$$T \left(\begin{array}{|c|c|c|c|c|} \hline c_{1,1} & c_{1,2} & c_{1,3} & c_{1,4} & c_{1,5} \\ \hline c_{2,1} & c_{2,2} & c_{2,3} & c_{2,4} & c_{2,5} \\ \hline c_{3,1} & c_{3,2} & c_{3,3} & & \\ \hline c_{4,1} & & & & \\ \hline c_{5,1} & & & & \\ \hline c_{6,1} & & & & \\ \hline \end{array} \right) = \begin{array}{|c|c|c|c|c|c|} \hline c_{1,2} & c_{1,3} & c_{1,4} & c_{1,5} & 0 & \\ \hline c_{2,2} & c_{2,3} & c_{2,4} & c_{2,5} & 0 & \\ \hline c_{3,2} & c_{3,3} & 0 & & & \\ \hline 0 & & & & & \\ \hline 0 & & & & & \\ \hline 0 & & & & & \\ \hline \end{array}. \quad (8.5)$$

It is immediate from this picture that the image of T consists of all F -linear combinations of the linearly independent basis vectors

$$\{x(1, 1), x(1, 2), x(1, 3), x(1, 4), x(2, 1), x(2, 2), x(2, 3), x(2, 4), x(3, 1), x(3, 2)\},$$

so that this set is a basis of $\text{img}(T)$. In general, we see pictorially or via (8.4) that the set

$$\{x(i, j - 1) : 1 \leq i \leq \ell(\mu), 2 \leq j \leq \mu_i\} = \{x(i, j) : (i, j) \in D(\mu) \text{ and } (i, j + 1) \in D(\mu)\}$$

is a basis for $\text{img}(T_{X, \mu})$. These are the basis vectors occupying all cells of $D(\mu)$ excluding the rightmost cell in each row.

What happens if we apply T again? In our example,

$$T^2 \left(\begin{array}{|c|c|c|c|c|} \hline c_{1,1} & c_{1,2} & c_{1,3} & c_{1,4} & c_{1,5} \\ \hline c_{2,1} & c_{2,2} & c_{2,3} & c_{2,4} & c_{2,5} \\ \hline c_{3,1} & c_{3,2} & c_{3,3} & & \\ \hline c_{4,1} & & & & \\ \hline c_{5,1} & & & & \\ \hline c_{6,1} & & & & \\ \hline \end{array} \right) = \begin{array}{|c|c|c|c|c|} \hline c_{1,3} & c_{1,4} & c_{1,5} & 0 & 0 \\ \hline c_{2,3} & c_{2,4} & c_{2,5} & 0 & 0 \\ \hline c_{3,3} & 0 & 0 & & \\ \hline 0 & & & & \\ \hline 0 & & & & \\ \hline 0 & & & & \\ \hline \end{array} .$$

So $\text{img}(T^2)$ consists of all F -linear combinations of $x(1, 1)$, $x(1, 2)$, $x(1, 3)$, $x(2, 1)$, $x(2, 2)$, $x(2, 3)$, and $x(3, 1)$, and these vectors form a basis for this subspace. Similarly, for any X and μ , a basis for $\text{img}(T_{X, \mu}^2)$ is the set of all $x(i, j)$ occupying cells $(i, j) \in D(\mu)$ that are not one of the *two* rightmost cells in each row.

More generally, iteration of (8.4) shows that for $k \in \mathbb{N}^+$,

$$\begin{aligned} T_{X, \mu}^k \left(\sum_{(i,j) \in D(\mu)} c(i, j)x(i, j) \right) &= \sum_{(i,j) \in D(\mu) : (i, j - k) \in D(\mu)} c(i, j)x(i, j - k) \\ &= \sum_{(i,j) \in D(\mu) : (i, j + k) \in D(\mu)} c(i, j + k)x(i, j) \quad (c(i, j) \in F), \end{aligned} \quad (8.6)$$

which shows that $\{x(i, j) : (i, j) \in D(\mu) \text{ and } (i, j + k) \in D(\mu)\}$ is an F -basis for $\text{img}(T_{X, \mu}^k)$. We get this basis by ignoring the k rightmost cells in each row of $D(\mu)$ and taking all remaining $x(i, j)$'s. In particular, taking $k = \mu_1 - 1$ and $k = \mu_1$ in (8.6) shows that $T_{X, \mu}$ is a nilpotent map of index μ_1 .

8.5 Computing Null Spaces via Partition Diagrams

Turning to null spaces, let us see how to use partition diagrams to find bases for $\text{Null}(T_{X, \mu}^k)$. For our example of $\mu = (5, 5, 3, 1, 1, 1)$, inspection of (8.5) reveals that $v = \sum_{(i,j) \in D(\mu)} c(i, j)x(i, j)$ is sent to zero by T iff

$$\begin{aligned} c(1, 2) &= c(1, 3) = c(1, 4) = c(1, 5) \\ &= c(2, 2) = c(2, 3) = c(2, 4) = c(2, 5) \\ &= c(3, 2) = c(3, 3) = 0. \end{aligned}$$

Pictorially,

$$\text{Null}(T) = \left\{ \begin{array}{|c|c|c|c|c|} \hline c_{1,1} & 0 & 0 & 0 & 0 \\ \hline c_{2,1} & 0 & 0 & 0 & 0 \\ \hline c_{3,1} & 0 & 0 & & \\ \hline c_{4,1} & & & & \\ \hline c_{5,1} & & & & \\ \hline c_{6,1} & & & & \\ \hline \end{array} : c(i, 1) \in F \right\}. \quad (8.7)$$

In general, $\text{Null}(T_{X,\mu})$ consists of all F -linear combinations of basis vectors residing in the leftmost column of $D(\mu)$, so that $\{x(i, 1) : 1 \leq i \leq \ell(\mu) = \mu'_1\}$ is a basis for $\text{Null}(T_{X,\mu})$.

Similarly, applying T^2 to v will produce zero iff $c(i, j) = 0$ for all $j > 2$, so that a basis of $\text{Null}(T^2)$ consists of all basis vectors in the first two columns of $D(\mu)$. In general, for all $k \geq 1$, the set

$$\{x(i, j) : (i, j) \in D(\mu) \text{ and } j \leq k\} \quad (8.8)$$

is a basis for $\text{Null}(T_{X,\mu}^k)$, and therefore

$$\dim(\text{Null}(T_{X,\mu}^k)) = \mu'_1 + \mu'_2 + \cdots + \mu'_k \quad (k \geq 1). \quad (8.9)$$

To prove this assertion without pictures, note from (8.6) that $T_{X,\mu}^k(v) = 0$ iff $c(i, j) = 0$ for all $(i, j) \in D(\mu)$ such that $(i, j - k) \in D(\mu)$, which holds iff v is a linear combination of the $x(i, j)$ with $(i, j) \in D(\mu)$ and $j \leq k$.

8.6 Classification of Nilpotent Maps (Stage 1)

We are now ready to start the proof of the **classification of nilpotent linear maps**. We want to show that *given any field F , any n -dimensional F -vector space V , and any nilpotent linear map $T : V \rightarrow V$, there exists a unique partition μ of n such that $[T]_X = J(0; \mu)$ for some ordered basis X of V* . As we saw at the end of §8.3, the conclusion is equivalent to the existence of a unique partition μ such that $T = T_{X,\mu}$ for some ordered basis X of V .

We can prove uniqueness of μ right away. Suppose $T : V \rightarrow V$ is nilpotent, and $T = T_{X,\mu} = T_{Y,\nu}$ for some partitions μ and ν and some ordered bases X and Y of V . For all $k \geq 1$, (8.9) shows that

$$\mu'_1 + \cdots + \mu'_k = \dim(\text{Null}(T^k)) = \nu'_1 + \cdots + \nu'_k.$$

Then $\mu = \nu$ follows from (8.3).

Proving existence of μ and X is more subtle. We will use induction on $n = \dim(V)$. For the base case $n = 0$, we must have $V = \{0\}$ and $T = 0$. Then $T = T_{X,\mu}$ holds if we choose $X = \emptyset$ and $\mu = ()$, which is a partition of zero of length zero. If the reader dislikes that base case, check that for $n = 1$ we will have $T = T_{X,\mu}$ where $X = \{x\}$ is any nonzero vector in V and $\mu = (1)$ is the unique partition of 1.

In the rest of the proof, we assume $n = \dim(V) > 1$ and that the existence assertion is already known to hold for all F -vector spaces of smaller dimension than n . Given a nilpotent linear map $T : V \rightarrow V$, we will build μ and X such that $T = T_{X,\mu}$ in three stages. In stage 1, we consider the subspace $V_1 = \text{img}(T) = T[V]$. Since $\dim(V) > 0$, we cannot have $V_1 = V$; for otherwise $T[V] = V$ would give $T^k[V] = V \neq \{0\}$ for all $k \in \mathbb{N}^+$, contradicting the nilpotence of T . Furthermore, T maps $V_1 = T[V]$ into itself, so we can

consider the restricted function $T|V_1 : V_1 \rightarrow V_1$, given by $T|V_1(v) = T(v)$ for all $v \in V_1$. The map $T|V_1$ is still linear and still nilpotent. Since V_1 is a proper subspace of V , we know $n_1 = \dim(V_1) < \dim(V)$, and the induction hypothesis gives us an ordered basis X_1 of V_1 and a partition $\nu = (\nu_1, \dots, \nu_s)$ of n_1 such that $T|V_1 = T_{X_1, \nu}$. For example, $D(\nu)$ and X_1 might look like this:

$$D(\nu) = \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array} \quad \begin{array}{ccccccccc} 0 & \xleftarrow{T|V_1} & x_{1,1} & \xleftarrow{T|V_1} & x_{1,2} & \xleftarrow{T|V_1} & x_{1,3} & \xleftarrow{T|V_1} & x_{1,4} \\ 0 & \xleftarrow{T|V_1} & x_{2,1} & \xleftarrow{T|V_1} & x_{2,2} & \xleftarrow{T|V_1} & x_{2,3} & \xleftarrow{T|V_1} & x_{2,4} \\ 0 & \xleftarrow{T|V_1} & x_{3,1} & \xleftarrow{T|V_1} & x_{3,2} & & & & \\ \hline \end{array}$$

8.7 Classification of Nilpotent Maps (Stage 2)

Continuing the proof, we are still trying to show that $T : V \rightarrow V$ has the form $T_{X, \mu}$ for some partition μ of n . If this conclusion were true, we know from §8.4 that $D(\nu)$ must consist of all cells of $D(\mu)$ excluding the rightmost cell of each row. The next stage of the proof is to recover the rows of the unknown partition μ of length 1, if there are any.

To do this, we study $V_2 = V_1 + \text{Null}(T) = \{u + v : u \in V_1, v \in \text{Null}(T)\}$, which is a subspace of V containing $V_1 = \text{img}(T)$. Write $n_2 = \dim(V_2)$. The subspace V_2 is spanned by $X_1 \cup \text{Null}(T)$, so we can extend the ordered basis X_1 for V_1 to an ordered basis X_2 for V_2 by adding appropriate vectors from $\text{Null}(T)$ to the end of the list X_1 . Suppose t additional vectors are needed, so $n_2 = n_1 + t$. Let ρ be the partition of n_2 obtained by adding t parts of size 1 to the end of ν . Note that $T|V_2$ is nilpotent and maps V_2 to itself (in fact, it maps V_2 into V_1). We can see that $T|V_2 = T_{X_2, \rho}$ since both sides are linear maps having the same effect on basis vectors in X_2 (in particular, both sides send the new vectors in $X_2 \sim X_1$ to zero). For example, $D(\rho)$ and X_2 might look like this (where new boxes are marked with stars):

$$D(\rho) = \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \star & & & \\ \hline \star & & & \\ \hline \star & & & \\ \hline \end{array} \quad \begin{array}{ccccccccc} 0 & \xleftarrow{T|V_2} & x_{1,1} & \xleftarrow{T|V_2} & x_{1,2} & \xleftarrow{T|V_2} & x_{1,3} & \xleftarrow{T|V_2} & x_{1,4} \\ 0 & \xleftarrow{T|V_2} & x_{2,1} & \xleftarrow{T|V_2} & x_{2,2} & \xleftarrow{T|V_2} & x_{2,3} & \xleftarrow{T|V_2} & x_{2,4} \\ 0 & \xleftarrow{T|V_2} & x_{3,1} & \xleftarrow{T|V_2} & x_{3,2} & & & & \\ 0 & \xleftarrow{T|V_2} & x_{4,1} & & & & & & \\ 0 & \xleftarrow{T|V_2} & x_{5,1} & & & & & & \\ 0 & \xleftarrow{T|V_2} & x_{6,1} & & & & & & \\ \hline \end{array}$$

We claim that $n_2 = n - s$, where $s = \ell(\nu)$ is the number of rows in $D(\nu)$. To prove the claim, recall that $V_2 = \text{img}(T) + \text{Null}(T)$, so (using Exercise 24)

$$n_2 = \dim(V_2) = \dim(\text{img}(T)) + \dim(\text{Null}(T)) - \dim(\text{img}(T) \cap \text{Null}(T)). \quad (8.10)$$

By the rank-nullity theorem (see §1.8), $\dim(\text{img}(T)) + \dim(\text{Null}(T)) = \dim(V) = n$. On the other hand, $\text{Null}(T|V_1) = V_1 \cap \text{Null}(T) = \text{img}(T) \cap \text{Null}(T)$. Applying (8.9) to the map $T|V_1 = T_{X_1, \nu}$ (with $k = 1$), we get

$$\dim(\text{img}(T) \cap \text{Null}(T)) = \dim(\text{Null}(T|V_1)) = \nu'_1 = \ell(\nu) = s.$$

Putting these formulas into (8.10) gives $n_2 = n - s$, as claimed.

8.8 Classification of Nilpotent Maps (Stage 3)

We are finally ready to find a partition μ of n and an ordered basis X of V for which $T = T_{X,\mu}$. Define μ by adding one to each of the first s parts of ρ , so $\mu = (\nu_1 + 1, \nu_2 + 1, \dots, \nu_s + 1, 1, \dots, 1)$ where there are t ones at the end. Extend the indexed basis $X_2 = \{x(i,j) : (i,j) \in D(\rho)\}$ to an indexed set $X = \{x(i,j) : (i,j) \in D(\mu)\}$ by letting $x(i,\mu_i)$ be any vector in V with $T(x(i,\mu_i)) = x(i,\mu_i - 1)$, for $1 \leq i \leq s$. Such vectors must exist, since $(i,\mu_i - 1) \in D(\nu)$ means that $x(i,\mu_i - 1) \in X_1 \subseteq \text{img}(T)$. For example, $D(\mu)$ and X might look like this (where new boxes are marked with stars):

$$D(\mu) = \begin{array}{|c|c|c|c|c|} \hline & & & & \star \\ \hline \end{array} \quad \begin{aligned} 0 &\xleftarrow{T} x_{1,1} \xleftarrow{T} x_{1,2} \xleftarrow{T} x_{1,3} \xleftarrow{T} x_{1,4} \xleftarrow{T} x_{1,5} \\ 0 &\xleftarrow{T} x_{2,1} \xleftarrow{T} x_{2,2} \xleftarrow{T} x_{2,3} \xleftarrow{T} x_{2,4} \xleftarrow{T} x_{2,5} \\ 0 &\xleftarrow{T} x_{3,1} \xleftarrow{T} x_{3,2} \xleftarrow{T} x_{3,3} \\ 0 &\xleftarrow{T} x_{4,1} \\ 0 &\xleftarrow{T} x_{5,1} \\ 0 &\xleftarrow{T} x_{6,1} \end{aligned}$$

If we can show that the indexed set X is a basis of V , then $T = T_{X,\mu}$ will follow since both linear maps agree on all basis vectors in X . By the claim in the last section, μ is a partition of $n_2+s = n = \dim(V)$. Hence, if we can show that the indexed set $X = \{x(i,j) : (i,j) \in D(\mu)\}$ is linearly independent, then the n vectors $x(i,j)$ must be distinct and form a basis of V .

To check independence, assume that for some scalars $c(i,j) \in F$,

$$0 = \sum_{(i,j) \in D(\mu)} c(i,j)x(i,j).$$

We must show every $c(i,j)$ is zero. Applying the linear map T gives

$$0 = \sum_{(i,j) \in D(\mu) : (i,j+1) \in D(\mu)} c(i,j+1)x(i,j) = \sum_{(i,j) \in D(\nu)} c(i,j+1)x(i,j)$$

(cf. (8.5) and (8.6)). Since $X_1 = \{x(i,j) : (i,j) \in D(\nu)\}$ is known to be linearly independent already (by induction), we see that $c(i,j+1) = 0$ for all $(i,j) \in D(\nu)$. Our original linear combination now reduces to

$$0 = \sum_{i=1}^{\ell(\mu)} c(i,1)x(i,1) = \sum_{(i,j) \in D(\rho) : j=1} c(i,j)x(i,j)$$

(cf. (8.7)). But $X_2 = \{x(i,j) : (i,j) \in D(\rho)\}$ is linearly independent by construction, so we deduce that $c(i,1) = 0$ for all i . Hence the indexed set X is linearly independent, and the existence proof is finally complete.

Let us mention one corollary of the classification of nilpotent linear maps. If V is n -dimensional and $T : V \rightarrow V$ is a nilpotent linear map, then $T^n = 0$, i.e., T has index of nilpotence at most n . To see why this holds, note $T = T_{X,\mu}$ for some ordered basis X and some partition μ of n . The diagram $D(\mu)$ can have at most n nonempty columns. By (8.9), $\text{Null}(T^n)$ has dimension $\mu'_1 + \dots + \mu'_n = n$, so $\text{Null}(T^n)$ must be all of V . This means that $T^n = 0$ as claimed.

8.9 Fitting's Lemma

We will prove the existence part of the Jordan canonical form theorem by combining the classification of nilpotent linear maps with a result called *Fitting's lemma*. Before stating this result, we must discuss some preliminary concepts. Let V be a vector space over a field F . Given two subspaces W and Z of V , the *sum* of these subspaces is the subspace $W + Z = \{w + z : w \in W, z \in Z\}$. We say this sum is a *direct sum*, denoted $W \oplus Z$, iff $W \cap Z = \{0_V\}$. In this case, one can check that if X_1 is an ordered basis of W and X_2 is an ordered basis of Z , then the list consisting of the vectors in X_1 followed by the vectors in X_2 is an ordered basis of $W \oplus Z$. In particular, $\dim(W \oplus Z) = \dim(W) + \dim(Z)$ when W and Z are finite-dimensional.

Given any linear map $S : V \rightarrow V$, we say a subspace W of V is *S -invariant* iff $S(w) \in W$ for all $w \in W$. In this situation, the restriction of S to W , denoted $S|W$, maps W into itself and is a linear map from W to W . Now assume $n = \dim(V)$ is finite. Recall that $\text{Null}(S) = \ker(S)$ and $\text{img}(S)$ are always subspaces of V , and by the rank-nullity theorem (see §1.8), $\dim(\text{Null}(S)) + \dim(\text{img}(S)) = n = \dim(V)$. If $\text{Null}(S) \cap \text{img}(S) = \{0\}$, then we would have a direct sum $\text{Null}(S) \oplus \text{img}(S)$ of dimension n , which must therefore be the entire space V . However, there is no guarantee that $\text{Null}(S)$ and $\text{img}(S)$ have zero intersection. For example, if S is the nilpotent map defined in (8.5), then $\text{Null}(S) \cap \text{img}(S)$ is a three-dimensional subspace spanned by $x(1, 1)$, $x(2, 1)$, and $x(3, 1)$.

We can rectify this situation by imposing an appropriate hypothesis on S . Specifically, if $S : V \rightarrow V$ is a linear map such that $\text{Null}(S) = \text{Null}(S^2)$, then $V = \text{Null}(S) \oplus \text{img}(S)$. By the comments in the preceding paragraph, it is enough to show that $\text{Null}(S) \cap \text{img}(S) = \{0\}$. Suppose $z \in \text{Null}(S) \cap \text{img}(S)$. On one hand, $S(z) = 0$. On the other hand, $z = S(y)$ for some $y \in V$. Now, $S^2(y) = S(S(y)) = S(z) = 0$ means that $y \in \text{Null}(S^2) = \text{Null}(S)$. Then $z = S(y) = 0$, so $\text{Null}(S) \cap \text{img}(S) = \{0\}$ as needed.

We now prove **Fitting's lemma**: given an n -dimensional F -vector space V and a linear map $U : V \rightarrow V$, there exist U -invariant subspaces Z and W of V such that $V = Z \oplus W$, $U|Z$ is nilpotent, and $U|W$ is an isomorphism.

If $\text{Null}(U) = \{0\}$, then $\dim(\text{img}(U)) = n$ by the rank-nullity theorem, so $\text{img}(U) = V$ and U is an isomorphism. So we may take $Z = \{0\}$ and $W = V$ in this situation. Now suppose $\text{Null}(U)$ is a nonzero subspace of V . Using the readily verified fact that $\text{Null}(U^k) \subseteq \text{Null}(U^{k+1})$ for all $k \geq 1$, we have a chain of subspaces

$$\{0\} \neq \text{Null}(U) \subseteq \text{Null}(U^2) \subseteq \text{Null}(U^3) \subseteq \cdots \subseteq \text{Null}(U^k) \subseteq \cdots \subseteq V.$$

All these subspaces are finite-dimensional, so there must exist $m \geq 1$ with $\text{Null}(U^k) = \text{Null}(U^m)$ for all $k \geq m$. In particular, taking $k = 2m$, we can apply our earlier result to the linear map $S = U^m$ to conclude that $V = \text{Null}(S) \oplus \text{img}(S) = \text{Null}(U^m) \oplus \text{img}(U^m)$. Let $Z = \text{Null}(U^m) \neq \{0\}$ and $W = \text{img}(U^m)$. Given $z \in Z$, we know $U^m(z) = 0$, so $U^m(U(z)) = U(U^m(z)) = U(0) = 0$, so $U(z) \in Z$. Given $w \in W$, we know $w = U^m(v)$ for some $v \in V$, so $U(w) = U(U^m(v)) = U^m(U(v)) \in W$. Thus Z and W are U -invariant subspaces. By the very definition of Z , $U|Z$ is a nilpotent linear map of index at most m . On the other hand, suppose $w \in W$ lies in $\text{Null}(U|W)$. Then $U(w) = 0$ implies $U^m(w) = 0$, so $w \in \text{Null}(U^m) \cap \text{img}(U^m) = \{0\}$. So $\ker(U|W) = \{0\}$, which means $U|W : W \rightarrow W$ is injective and hence surjective. Thus, $U|W$ is an isomorphism.

8.10 Existence of Jordan Canonical Forms

In this section, we prove the existence assertion in the Jordan canonical form theorem (stated in the introduction to this chapter). Before doing so, we must recall some basic facts about eigenvalues. Given a field F and a matrix $A \in M_n(F)$, a scalar $c \in F$ is called an *eigenvalue* of A iff $Av = cv$ for some nonzero $n \times 1$ column vector v . Any such v is called an *eigenvector* of A associated with the eigenvalue c . The eigenvalues of A are precisely the roots of the characteristic polynomial $\chi_A = \det(xI_n - A) \in F[x]$ that lie in the field F . This polynomial has degree n , so A has at most n distinct eigenvalues in F . For the field $F = \mathbb{C}$ of complex numbers, the polynomial χ_A always has a complex root by the fundamental theorem of algebra. So every $n \times n$ complex matrix A has between 1 and n complex eigenvalues. For $A \in M_n(F)$, let $\text{Spec}(A)$ (the *spectrum* of A) be the set of all eigenvalues of A . When A is triangular, $\text{Spec}(A)$ is the set of scalars appearing on the main diagonal of A .

Next, if V is an n -dimensional F -vector space with $n > 0$ and $T : V \rightarrow V$ is a linear map, an *eigenvalue* of T is a scalar $c \in F$ such that $T(v) = cv$ for some nonzero $v \in V$; any such v is called an *eigenvector* of T associated with the eigenvalue c . If X is any ordered basis of V and $A = [T]_X$ is the matrix of T relative to X , then T and A have the same eigenvalues. This follows since $T(v) = cv$ iff $A[v]_X = c[v]_X$, where $[v]_X$ is the column vector giving the coordinates of v relative to X . In particular, every linear map T on an n -dimensional complex vector space has between 1 and n complex eigenvalues. Let $\text{Spec}(T)$ (the *spectrum* of T) be the set of all eigenvalues of T .

We now prove: *for every finite-dimensional complex vector space V and every linear map $T : V \rightarrow V$, there exists an ordered basis X of V such that $[T]_X$ is a Jordan canonical form.* The proof is by induction on $n = \dim(V)$. For $n = 0$ and $n = 1$, the result follows immediately. Now assume $n > 1$ and the result is already known for all complex vector spaces of smaller dimension than n .

Pick a fixed eigenvalue c of the given linear map T , and let $v \neq 0$ be an associated eigenvector. Let $U = T - c\text{Id}_V$, where Id_V denotes the identity map on V . Applying Fitting's lemma to the linear map U , we get a direct sum $V = Z \oplus W$ where Z and W are U -invariant (hence also T -invariant) subspaces of V such that $U|Z$ is nilpotent and $U|W$ is an isomorphism. On one hand, $U|Z$ is nilpotent, so we know there is an ordered basis X_1 of Z and a partition μ of the integer $k = \dim(Z)$ such that $[U|Z]_{X_1} = J(0; \mu)$. It follows that

$$[T|Z]_{X_1} = [U|Z + c\text{Id}_Z]_{X_1} = [U|Z]_{X_1} + [c\text{Id}_Z]_{X_1} = J(0; \mu) + cI_k = J(c; \mu).$$

On the other hand, v cannot lie in W , since otherwise $U|W(v) = U(v) = 0 = U|W(0)$ contradicts the fact that $U|W$ is an isomorphism. So $W \neq V$, forcing $\dim(W) < \dim(V)$. By the induction hypothesis, there exists an ordered basis X_2 of W such that $[U|W]_{X_2}$ is a Jordan canonical form matrix J_1 . By the same calculation used above, we see that $[T|W]_{X_2}$ is the matrix J_2 obtained from J_1 by adding c to every diagonal entry. This new matrix is also a Jordan canonical form. Finally, taking X to be the concatenation of X_1 and X_2 , we know X is an ordered basis of V such that $[T]_X = \text{blk-diag}(J(c; \mu), J_2)$. This matrix is a Jordan canonical form, so the induction proof is complete.

The existence of Jordan canonical forms for linear maps implies the existence of Jordan canonical forms for matrices, as follows. Given $A \in M_n(\mathbb{C})$, let $T : \mathbb{C}^n \rightarrow \mathbb{C}^n$ be the linear map defined by $T(v) = Av$ for all column vectors $v \in \mathbb{C}^n$. Choose an ordered basis X of \mathbb{C}^n such that $J = [T]_X$ is a Jordan canonical form. We know $J = P^{-1}AP$ for some invertible $P \in M_n(\mathbb{C})$, so A is similar to a Jordan canonical form.

We should also point out that the only special feature of the field \mathbb{C} needed in this

proof was that every linear map on an n -dimensional \mathbb{C} -vector space (with $n > 0$) has an eigenvalue in \mathbb{C} . This follows from the fact that all non-constant polynomials in $\mathbb{C}[x]$ split into products of linear factors. The Jordan canonical form theorem extends to any field F having the latter property (such fields are called *algebraically closed*).

8.11 Uniqueness of Jordan Canonical Forms

Recall the uniqueness assertion in the Jordan canonical form theorem for linear maps: *if V is an n -dimensional complex vector space, $T : V \rightarrow V$ is linear, and X and Y are ordered bases of V such that $A = [T]_X$ and $B = [T]_Y$ are both Jordan canonical forms, then B is obtained from A by rearranging the Jordan blocks in A .*

Before proving this statement, we consider an example that conveys the idea of the proof. Suppose $X = (x_1, \dots, x_{14})$ and $A = [T]_X$ is the Jordan canonical form

$$\text{blk-diag}(J(7; 2), J(7; 2), J(7; 1), J(-4; 5), J(i; 3), J(i; 1)).$$

This matrix is triangular, so we can see that $\text{Spec}(T) = \text{Spec}(A) = \{7, -4, i\}$. Some eigenvectors of A and T associated with the eigenvalue 7 are x_1 , x_3 , and x_5 . In fact, one can check that the set of all eigenvectors for this eigenvalue is the set of all nonzero \mathbb{C} -linear combinations of x_1 , x_3 , and x_5 . On the other hand, x_2 and x_4 are not eigenvectors for the eigenvalue 7, even though these basis vectors “belong to” Jordan blocks with 7 on the diagonal.

The key to the uniqueness proof is realizing that the subspace spanned by x_1, \dots, x_5 can be described in terms of just the linear map T , not the ordered basis X or the matrix A . To see how this is done, consider the linear map $U = T - 7 \text{Id}_V$. The matrix of U relative to X is

$$C = [U]_X = A - 7I_{14} = \text{blk-diag}(J(0; 2), J(0; 2), J(0; 1), J(-11; 5), J(-7+i; 3), J(-7+i; 1)).$$

We can also write $C = \text{blk-diag}(J(0; \mu), C_1)$, where $\mu = (2, 2, 1)$ and C_1 is a triangular matrix with all diagonal entries nonzero. What is the null space of C^{14} ? We compute $C^{14} = \text{blk-diag}(0_{5 \times 5}, C_2)$, where $C_2 = C_1^{14}$ is still a triangular matrix with all diagonal entries nonzero. One can now check that $\text{Null}(C^{14})$ consists of all column vectors $v \in \mathbb{C}^{14}$ with v_1, \dots, v_5 arbitrary and $v_6 = v_7 = \dots = v_{14} = 0$. Translating back to the linear map T , this means that $\text{Null}((T - 7 \text{Id}_V)^{14})$ is the subspace of V with basis $(x_1, x_2, x_3, x_4, x_5)$. Since the restriction of $T - 7 \text{Id}_V$ to this subspace is nilpotent, we can appeal to the known uniqueness result for nilpotent maps to see that the partition $\mu = (2, 2, 1)$ is unique.

Let us now turn to the general uniqueness proof. Assume the setup in the first paragraph of this section. First of all, $\text{Spec}(A) = \text{Spec}(T) = \text{Spec}(B)$, so that A and B have the same set of diagonal entries (ignoring multiplicities), namely the set of eigenvalues of the map T .

Let $c \in \text{Spec}(T)$ be any fixed eigenvalue. Write $U = T - c \text{Id}_V$. Solely for notational convenience, we can reorder the ordered basis X and the ordered basis Y so that all the Jordan blocks for c in A and B occur first, with the block sizes weakly decreasing, say $A = \text{blk-diag}(J(c; \mu), A')$ for some partition μ of some $k \in \mathbb{N}^+$ and $B = \text{blk-diag}(J(c; \nu), B')$ for some partition ν of some $m \in \mathbb{N}^+$. Since c is fixed but arbitrary, uniqueness will be proved if we can show $k = m$ and $\mu = \nu$. Writing the reordered bases as $X = (x_1, \dots, x_n)$ and $Y = (y_1, \dots, y_n)$, we claim that $X_1 = (x_1, \dots, x_k)$ and $Y_1 = (y_1, \dots, y_m)$ are both bases for the same subspace $Z = \text{Null}(U^n)$ of V . Assuming this claim is true, it follows that $k = m$, $[U|Z]_{X_1} = J(0; \mu)$, and $[U|Z]_{Y_1} = J(0; \nu)$. Since $U|Z$ is a nilpotent linear map,

our previously proved uniqueness result for nilpotent operators allows us to conclude that $\mu = \nu$.

So we need only prove the claim that $X_1 = (x_1, \dots, x_k)$ is a basis for $Z = \text{Null}(U^n)$ (the corresponding assertion for Y_1 is proved in the same way). Since $U = T - c\text{Id}_V$, $[U]_X = \text{blk-diag}(J(0; \mu), A' - cI_{n-k})$ where $A' - cI_{n-k}$ is an upper-triangular matrix with no zeroes on its diagonal (because all Jordan blocks for c occur in the first k rows of A). Raising this $n \times n$ matrix to the power n , we get $[U^n]_X = [U]_X^n = \text{blk-diag}(0_{k \times k}, A'')$, where $A'' \in M_{n-k}(\mathbb{C})$ is triangular with nonzero entries on its diagonal. From the form of the matrix $[U^n]_X$, we see immediately that U^n sends x_1, \dots, x_k , and every linear combination of these vectors to zero. So the span of X_1 is contained in $\text{Null}(U^n) = Z$. For the reverse inclusion, suppose $v = \sum_{i=1}^n c_i x_i$ (where $c_i \in \mathbb{C}$) is not in the span of X_1 . This means there exists $s > k$ with $c_s \neq 0$; choose a maximal index s with this property. Since $[U^n]_X$ is triangular with a nonzero entry in the s, s -position, we see that (for some scalars d_i) $U^n(v) = \sum_{i < s} d_i x_i + c_s A''(s, s)x_s$ is not zero, hence $v \notin Z$. We now know Z is the subspace spanned by the list X_1 . Since this list is linearly independent (being a sublist of the ordered basis X), X_1 is an ordered basis of Z .

The uniqueness of Jordan canonical forms for linear maps implies the uniqueness result for matrices, as follows. Given $A \in M_n(\mathbb{C})$, let $T : \mathbb{C}^n \rightarrow \mathbb{C}^n$ be the linear map $T(v) = Av$ for $v \in \mathbb{C}^n$. The set of matrices similar to A is exactly the set of matrices $[T]_X$ as X ranges over all ordered bases of \mathbb{C}^n (see Chapter 6). So all Jordan canonical forms in the similarity class of A represent the map T relative to appropriate bases, and these forms therefore differ from one another merely by rearranging the Jordan blocks.

8.12 Computing Jordan Canonical Forms

In this section, we briefly discuss how one might actually compute a Jordan canonical form of a specific matrix $A \in M_n(\mathbb{C})$. More specifically, we want to find a Jordan canonical form $\mathbf{J} \in M_n(\mathbb{C})$ and an invertible matrix P with $\mathbf{J} = P^{-1}AP$. If all we need is the matrix \mathbf{J} , we can proceed as follows. Examining the uniqueness proof given above, we observe first that the diagonal entries of \mathbf{J} must be the eigenvalues of A . These eigenvalues can be found, in principle, by computing the roots of the characteristic polynomial of A . In reality, for large n , it may be difficult or impossible to find the exact eigenvalues of A . But for the present discussion, let us assume that $\text{Spec}(A)$ can be found.

For each $c \in \text{Spec}(A)$, we need to find the Jordan blocks in \mathbf{J} of the form $J(c; m)$. Assuming these blocks occur in decreasing order of size, they collectively constitute a Jordan matrix of the form $J(c; \mu)$ for some unknown partition μ depending on c . Looking at the uniqueness proof again, we are led to consider the null space Z of the matrix $(A - cI_n)^n$. The proof shows that the nilpotent linear map $U : Z \rightarrow Z$ defined by $U(z) = Az - cz$ for $z \in Z$ has matrix $J(0; \mu)$ relative to some ordered basis of Z . We can find μ by using (8.9), which states that $\mu'_1 + \dots + \mu'_k = \dim(\text{Null}(U^k))$ for all $k \geq 1$. Translating back to matrices, we need only compute bases for all the null spaces $\text{Null}((A - cI_n)^k)$ for $1 \leq k \leq n$, which can be done by Gaussian elimination. Letting d_k be the dimension of the k 'th null space (with $d_0 = 0$), we compute $\mu'_k = d_k - d_{k-1}$ for all $k \geq 1$, from which $D(\mu)$ and hence μ are readily found.

Finding the transition matrix P requires more work. Since \mathbf{J} is now known, one straightforward but inefficient approach is to treat the entries of P as unknowns, which can be found by solving the linear system $P\mathbf{J} = AP$. Another method first finds bases for the various null spaces $\text{Null}((A - cI_n)^n)$. Using these bases to fill the columns of a matrix

Q , we will have $Q^{-1}AQ = \text{blk-diag}(A_1, \dots, A_k)$, where each A_i has a single eigenvalue c_i and $A_i - c_i I$ is nilpotent. Considering each block separately, we are reduced to solving the following problem: given a nilpotent matrix $B \in M_m(\mathbb{C})$, find a specific invertible matrix R and a partition μ of m such that $R^{-1}BR = J(0; \mu)$.

This problem can be solved by a recursive algorithm that implements the three-stage inductive proof of the classification theorem for nilpotent linear operators. First, we recursively find an ordered basis X_1 for $V_1 = \text{img}(B)$ and a partition $\nu = (\nu_1, \dots, \nu_s)$ such that the map $(x \mapsto Bx : x \in V_1)$ has matrix $J(0; \nu)$ relative to the basis X_1 . Second, we compute a basis for $\text{Null}(B)$ and use appropriate vectors from this basis to augment X_1 to a basis X_2 for $V_2 = \text{img}(B) + \text{Null}(B)$. If t new basis vectors are added, we define ρ by adding t new parts of size 1 to the end of ν . Third, writing $X_2 = \{x(i, j) : (i, j) \in D(\rho)\}$, we solve linear equations to find vectors $x(i, \nu_j + 1)$ such that $Bx(i, \nu_j + 1) = x(i, \nu_j)$ for $1 \leq i \leq s$. We define $\mu = (\nu_1 + 1, \dots, \nu_s + 1, 1, \dots, 1)$ (where there are t parts of size one) and $X = \{x(i, j) : (i, j) \in D(\mu)\}$. Finally, we obtain R by placing the column vectors $x(i, j)$ into a matrix, starting with the first row of $D(\mu)$ and working left to right, top to bottom. The reader can check that every step in this paragraph can be implemented by solving an explicitly computable system of linear equations. Admittedly, the entire algorithm requires a very substantial amount of computation.

For example, consider the matrix

$$A = \begin{bmatrix} -2 & -7 & 2 & 1 & 7 & -8 \\ -3.5 & -1.5 & -0.5 & 0.5 & 5 & 0 \\ -1 & -4 & 1 & 1 & 4 & -5 \\ 6.5 & 18.5 & -10.5 & -0.5 & -14 & 31 \\ -0.5 & -0.5 & -3.5 & 0.5 & 4 & 3 \\ 3 & 3 & -3 & 0 & -3 & 5 \end{bmatrix}. \quad (8.11)$$

Using a computer algebra system, we compute $\chi_A = (x + 1)^2(x - 2)^4$, so $\text{Spec}(A) = \{-1, 2\}$. Row reduction of $(A + I)^6$ shows that the null space of this matrix consists of column vectors

$$\{(s, s - t, s + t, -s - 2t, s, t) : s, t \in \mathbb{C}\},$$

so $v_1 = (1, 1, 1, -1, 1, 0)$ and $v_2 = (0, -1, 1, -2, 0, 1)$ form a basis for this null space. Similarly, row reduction of $(A - 2I)^6$ yields a null space

$$\{(a - b/4 + c/4 + 3d/4, b/4 + 3c/4 - 7d/4, a, b, c, d) : a, b, c, d \in \mathbb{C}\},$$

so

$$v_3 = (1, 0, 1, 0, 0, 0), \quad v_4 = (-1, 1, 0, 4, 0, 0), \quad v_5 = (1, 3, 0, 0, 4, 0), \quad v_6 = (3, -7, 0, 0, 0, 4)$$

form a basis for this null space. Letting P_1 be the matrix with columns v_1, \dots, v_6 , we find

$$A_1 = P_1^{-1}AP_1 = \begin{bmatrix} -1 & -1 & 0 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 3 & 5 \\ 0 & 0 & -1 & 2.5 & 1.5 & 3.5 \\ 0 & 0 & -1 & 0.5 & 3.5 & 3.5 \\ 0 & 0 & 0 & 0 & 0 & 2 \end{bmatrix}.$$

By inspection, replacing v_2 by $-v_2$ will convert the upper 2×2 block to $J(-1; 2)$. To deal with the lower 4×4 block (corresponding to the eigenvalue 2 of A), we consider the nilpotent matrix

$$B = \begin{bmatrix} -2 & 1 & 3 & 5 \\ -1 & 0.5 & 1.5 & 3.5 \\ -1 & 0.5 & 1.5 & 3.5 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

A recursive call to the algorithm produces a basis $x(1,1) = (2, 1, 1, 0)$, $x(1,2) = (2.5, 1.75, 1.75, 0)$ for $\text{img}(B)$ and associated partition $\nu = (2)$. Row-reduction of B shows that $\text{Null}(B)$ has a basis $((1, 2, 0, 0), (3, 0, 2, 0))$. Since the first basis vector is not a linear combination of $x(1,1)$ and $x(1,2)$, we set $x(2,1) = (1, 2, 0, 0)$ and $\rho = (2, 1)$. Finally, we solve $Bv = x(1,2)$ to obtain $x(1,3) = (2, 1, 1, 0.5)$. Now, taking

$$P_2 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 2.5 & 2 & 1 \\ 0 & 0 & 1 & 1.75 & 1 & 2 \\ 0 & 0 & 1 & 1.75 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0.5 & 0 \end{bmatrix}, \quad P = P_1 P_2 = \begin{bmatrix} 1 & 0 & 2 & 2.5 & 3.5 & -1 \\ 1 & 1 & 4 & 7 & 0.5 & 2 \\ 1 & -1 & 2 & 2.5 & 2 & 1 \\ -1 & 2 & 4 & 7 & 4 & 8 \\ 1 & 0 & 4 & 7 & 4 & 0 \\ 0 & -1 & 0 & 0 & 2 & 0 \end{bmatrix},$$

we find that $P^{-1}AP = \text{blk-diag}(J(-1; 2), J(2; 3), J(2; 1))$.

8.13 Application to Differential Equations

This section describes an application of the Jordan canonical form theorem to the solution of homogeneous systems of linear ordinary differential equations. We let x denote a column vector of unknown continuously differentiable functions $x_1, \dots, x_n : \mathbb{R} \rightarrow \mathbb{C}$. Given a fixed matrix $A \in M_n(\mathbb{C})$, we wish to find all solutions to $x' = Ax$, which encodes the system of differential equations

$$x'_i(t) = \sum_{j=1}^n A(i,j)x_j(t) \quad (1 \leq i \leq n, t \in \mathbb{R}).$$

This system could be solved easily if A were a diagonal matrix; for then $x'_i(t) = A(i,i)x_i(t)$ has general solution $x_i(t) = b_i e^{A(i,i)t}$, where b_i is any constant. More generally, if A is a Jordan block $J(c; n)$ for some $c \in \mathbb{C}$, we can “backsolve” for $x_n(t), x_{n-1}(t), \dots, x_1(t)$ as follows. First, $x'_n(t) = cx_n(t)$ implies $x_n(t) = b_n e^{ct}$ for some constant b_n . Second, $x'_{n-1}(t) = cx_{n-1}(t) + x_n(t) = cx_{n-1}(t) + b_n e^{ct}$ has general solution $x_{n-1}(t) = b_n t e^{ct} + b_{n-1} e^{ct}$ where b_{n-1} is any constant. Third, $x'_{n-2}(t) = cx_{n-2}(t) + x_{n-1}(t)$ has general solution $x_{n-2}(t) = (b_n/2)t^2 e^{ct} + b_{n-1} t e^{ct} + b_{n-2} e^{ct}$ where b_{n-2} is any constant. Continuing backwards, one can check by induction that

$$x_{n-k}(t) = \sum_{i=0}^k \frac{b_{n-i}}{(k-i)!} t^{k-i} e^{ct} \quad (0 \leq k < n), \tag{8.12}$$

where $b_n, \dots, b_1 \in \mathbb{C}$ are arbitrary constants.

Now, if $A = \text{blk-diag}(J(c_1; n_1), \dots, J(c_k; n_k))$, then one can use (8.12) to solve for the first n_1 functions x_1, \dots, x_{n_1} . The same formula applies to recover the functions $x_{n_1+1}, \dots, x_{n_1+n_2}$, starting with $x_{n_1+n_2}$ and working backwards, and similarly for all later blocks.

Finally, consider the case where $A \in M_n(\mathbb{C})$ is an arbitrary matrix. Using the algorithms in the previous section, we can find a Jordan canonical form \mathbf{J} and an invertible matrix P in $M_n(\mathbb{C})$ with $\mathbf{J} = P^{-1}AP$. Introduce a new column vector y of unknown functions y_1, \dots, y_n by the linear change of variable $y = P^{-1}x$, $x = Py$. By linearity of the derivative, $x' = Py'$, and we see that x solves $x' = Ax$ iff $Py' = A(Py)$ iff $y' = (P^{-1}AP)y$ iff y solves $y' = \mathbf{J}y$.

So, once we compute P and \mathbf{J} , y is given by the formulas above, and then $x = Py$ is the solution to the original system. We hasten to remark that there are usually more efficient ways of solving $x' = Ax$, especially when A has special structure, but a full discussion of this point is beyond the scope of this text.

For example, let us solve $x' = Ax$, where A is the matrix in (8.11). We computed a Jordan canonical form of A to be blk-diag($J(-1; 2), J(2; 3), J(2; 1)$). Let $y = P^{-1}x$, where P is the matrix found at the end of §8.12. Then $y' = \mathbf{J}y$ has general solution

$$\begin{aligned} y_1(t) &= c_1te^{-t} + c_2e^{-t}, & y_2(t) &= c_1e^{-t}, \\ y_3(t) &= (c_3/2)t^2e^{2t} + c_4te^{2t} + c_5e^{2t}, & y_4(t) &= c_3te^{2t} + c_4e^{2t}, & y_5(t) &= c_3e^{2t}, \\ y_6(t) &= c_6e^{2t}. \end{aligned}$$

where $c_1, \dots, c_6 \in \mathbb{C}$ are arbitrary constants, and the original system has solution $x = Py$.

8.14 Minimal Polynomials

Let V be an n -dimensional vector space over any field F , and let $T : V \rightarrow V$ be a linear map. Recall that T lies in the F -algebra $L(V)$ of all F -linear maps from V to V , which is finite-dimensional. Therefore, we know from §3.19 that there exists a unique monic polynomial $m_T \in F[x]$ of least degree such that $m_T(T) = 0$, and m_T divides all polynomials $g \in F[x]$ such that $g(T) = 0$. We can find m_T by searching the list of powers $(\text{Id}_V, T, T^2, T^3, \dots)$ for the lowest power T^k that is a linear combination of preceding powers of T . Similarly, any matrix $A \in M_n(F)$ has a minimal polynomial m_A . For any fixed ordered basis X of V , the map sending $T \in L(V)$ to $[T]_X \in M_n(F)$ is an algebra isomorphism (see Chapter 6). It follows that the linear map T has the same minimal polynomial as the matrix $[T]_X$, for any choice of X . By letting X vary over all ordered bases of V , we conclude from this that *similar matrices have the same minimal polynomial*.

Suppose A is a block-diagonal matrix blk-diag(A_1, \dots, A_s). For any $f \in F[x]$, one sees that $f(A) = \text{blk-diag}(f(A_1), \dots, f(A_s))$. Therefore, $f(A) = 0$ iff every $f(A_i) = 0$ iff m_{A_i} divides f for $1 \leq i \leq s$. It follows that $m_A = \text{lcm}(m_{A_1}, \dots, m_{A_s})$. In particular, suppose A is a diagonal matrix. Then we can take each A_i to be a 1×1 matrix. It is immediate that the minimal polynomial of the matrix $[c]$ is $x - c$, for any $c \in F$. Hence, for diagonal A , $m_A = \text{lcm}_{1 \leq i \leq n}(x - A(i, i))$. Since the diagonal entries of A are the eigenvalues of A , we can also write this as $m_A = \prod_{c \in \text{Spec}(A)} (x - c)$. More generally, any *diagonalizable* matrix A is similar to a diagonal matrix with the eigenvalues in $\text{Spec}(A)$ appearing on the main diagonal (each with a certain multiplicity). It follows that

$$\text{for diagonalizable } A \in M_n(F), m_A = \prod_{c \in \text{Spec}(A)} (x - c).$$

Note that this is a polynomial that splits into distinct linear factors in $F[x]$.

In fact, *for every matrix $A \in M_n(F)$, A is diagonalizable in $M_n(F)$ iff m_A splits into distinct linear factors in $F[x]$* . We use Jordan canonical forms to prove the backward direction for $F = \mathbb{C}$; a different proof is needed for general fields F (Exercise 48). Assume $A \in M_n(\mathbb{C})$ is not diagonalizable. We know A is similar to a Jordan canonical form $\mathbf{J} = \text{blk-diag}(J(c_1; n_1), \dots, J(c_s; n_s))$ that has the same minimal polynomial as A . The matrix \mathbf{J} cannot be diagonal, so some $n_i > 1$. We know $m_{\mathbf{J}} = \text{lcm}_{1 \leq i \leq s} m_{J(c_i; n_i)}$. So it will suffice to show that $m_{J(c; k)} = (x - c)^k$.

Evaluating the polynomial $g = (x - c)^k$ at $x = J(c; k)$ gives

$$g(J(c; k)) = (J(c; k) - cI_k)^k = J(0; k)^k = 0.$$

(The last equality can be seen by a matrix calculation, or by noting $J(0; k)$ is the matrix of the linear map $T_{X, (k)}$, which is nilpotent of index k .) So $m_{J(c; k)}$ must divide g . By unique factorization in $\mathbb{C}[x]$, the monic divisors of g in $\mathbb{C}[x]$ are $(x - c)^i$ where $0 \leq i \leq k$. Since $J(0; k)^i \neq 0$ for $i < k$, we see that g itself must be the minimal polynomial of $J(c; k)$.

Observe that the proof tells us how to compute m_A from a Jordan canonical form of A : $m_A = \prod_{c \in \text{Spec}(A)} (x - c)^{k(c)}$, where $k(c)$ is the maximum size of any Jordan block $J(c; k)$ appearing in a Jordan canonical form similar to A . On the other hand, consider the characteristic polynomial $\chi_A = \det(xI_n - A)$. Similar matrices have the same characteristic polynomial, so $\chi_A = \chi_J$. The characteristic polynomial of a Jordan block $J(c; k)$ is $(x - c)^k$, since $xI_n - J(c; k)$ is a $k \times k$ triangular matrix with all diagonal entries equal to $x - c$. For any block-diagonal matrix $B = \text{blk-diag}(B_1, \dots, B_s)$, computing the determinant shows that $\chi_B = \prod_{i=1}^s \chi_{B_i}$. Taking $B = \mathbf{J}$ here, we see that $\chi_A = \prod_{i=1}^s (x - c_i)^{n_i}$. Comparing to the earlier formula for m_A , we deduce the *Cayley–Hamilton theorem for complex matrices*: *for $A \in M_n(\mathbb{C})$, the minimal polynomial m_A divides the characteristic polynomial χ_A in $\mathbb{C}[x]$* .

Analogous results hold for linear maps $T \in L(V)$; in particular, *T is diagonalizable iff m_T splits into distinct linear factors in $F[x]$, in which case $m_T = \prod_{c \in \text{Spec}(T)} (x - c)$* . We use this theorem to prove that *if $T \in L(V)$ is diagonalizable and W is a T -invariant subspace of V , then $T|W$ is diagonalizable*. To prove this, write $g = m_T$ and $f = m_{T|W}$ in $F[x]$. We know $g(T)$ is the zero operator on V . Restricting to W , $g(T|W) = g(T)|W$ is the zero operator on W . So f , the minimal polynomial of $T|W$, must divide g in $F[x]$. By diagonalizability of T , g splits into a product of distinct linear factors in $F[x]$. Since f divides g , the unique factorization of f in $F[x]$ must also be a product of distinct linear factors (which are a subset of the factors for g). By the theorem, $T|W$ is diagonalizable.

8.15 Jordan–Chevalley Decomposition of a Linear Operator

This section proves an abstract version of the Jordan canonical form theorem that is needed in the theory of Lie algebras. Let V be an m -dimensional complex vector space. We will prove: *for each linear map $T : V \rightarrow V$, there exist unique linear maps T_d and T_n on V such that T_d is diagonalizable, T_n is nilpotent, $T = T_d + T_n$, and $T_d \circ T_n = T_n \circ T_d$. T_d is called the *diagonalizable part* of T , T_n is called the *nilpotent part* of T , and $T = T_d + T_n$ is called the *Jordan–Chevalley decomposition* of T* .

Fix the linear map T on V . To prove existence of T_d and T_n with the stated properties, write $\text{Spec}(T) = \{c_1, \dots, c_s\}$. By the Jordan canonical form theorem, we can pick an ordered basis X of V such that $[T]_X = \text{blk-diag}(J(c_1; \mu^{(1)}), J(c_2; \mu^{(2)}), \dots, J(c_s; \mu^{(s)}))$, where each $\mu^{(i)}$ is a partition of some $m_i \in \mathbb{N}^+$. Let X_1 consist of the first m_1 vectors in the list X , let X_2 consist of the next m_2 vectors in X , etc., so that X_s consists of the last m_s vectors in X . Define a linear map T_d on the basis X by letting $T_d(z_i) = c_i z_i$ for all $z_i \in X_i$, and extending by linearity. Define $T_n = T - T_d$, which is linear since it is a linear combination of two linear maps. Evidently $T = T_d + T_n$. Moreover, $[T_d]_X$ is the diagonal matrix $\text{blk-diag}(c_1 I_{m_1}, c_2 I_{m_2}, \dots, c_s I_{m_s})$, so T_d is diagonalizable. On the other hand,

$$[T_n]_X = [T]_X - [T_d]_X = \text{blk-diag}(J(0; \mu^{(1)}), J(0; \mu^{(2)}), \dots, J(0; \mu^{(s)})).$$

It follows that every basis vector in X gets sent to zero after applying T_n at most $\max(m_1, \dots, m_s)$ times, so that T_n is nilpotent. Finally, for each i , the scalar multiple of the identity $c_i I_{m_i}$ commutes with $J(0; \mu^{(i)})$. So the matrices $[T_d]_X$ and $[T_n]_X$ commute since each diagonal block of the first matrix commutes with the corresponding block of the second matrix. It follows that the linear maps T_d and T_n commute, completing the existence proof.

Turning to the uniqueness proof, let $T = T'_d + T'_n$ be any decomposition of T into the sum of a diagonalizable linear map T'_d and a nilpotent linear map T'_n that commute with each other. We must prove $T'_d = T_d$ and $T'_n = T_n$. Keep the notation of the previous paragraph. We saw in the uniqueness proof in §8.11 that each X_i is a basis of the T -invariant subspace $Z_i = \text{Null}((T - c_i \text{Id}_V)^m)$. Now, T'_d commutes with T , since

$$T'_d \circ T = T'_d \circ (T'_d + T'_n) = T'_d \circ T'_d + T'_d \circ T'_n = T'_d \circ T'_d + T'_n \circ T'_d = (T'_d + T'_n) \circ T'_d = T \circ T'_d.$$

Then T'_d also commutes with $T - c_i \text{Id}_V$ and any power of this linear map. It follows that each Z_i is a T'_d -invariant subspace, since $z \in Z_i$ implies $(T - c_i \text{Id}_V)^m(T'_d(z)) = T'_d((T - c_i \text{Id}_V)^m(z)) = T'_d(0) = 0$. Similarly, T'_n commutes with T , and so each Z_i is a T'_n -invariant subspace.

It now follows that $[T'_d]_X = \text{blk-diag}(B_1, \dots, B_s)$ and $[T'_n]_X = \text{blk-diag}(C_1, \dots, C_s)$ for certain matrices $B_i, C_i \in M_{m_i}(\mathbb{C})$. Now, B_i certainly commutes with $c_i I_{m_i}$ for all i , which says that T'_d commutes with T_d . Since T'_d also commutes with T , T'_d commutes with $T_n = T - T_d$. Similarly, T'_n commutes with T_d and T_n . Now, $T'_d + T'_n = T = T_d + T_n$ gives $T'_d - T_d = T_n - T'_n$. The left side of this equation is a diagonalizable linear map, since it is the difference of two commuting diagonalizable linear maps (Exercise 50). The right side of this equation is a nilpotent linear map, since it is the difference of two commuting nilpotent linear maps (Exercise 51). The two sides of the equation are equal, so we are considering a diagonalizable nilpotent linear map. The only such map is zero (as we saw in the introduction to this chapter), so $T'_d - T_d = 0 = T_n - T'_n$. Thus, $T'_d = T_d$ and $T_n = T'_n$, proving uniqueness.

The result just proved leads to the following Jordan–Chevalley decomposition theorem for matrices: *for any matrix $A \in M_n(\mathbb{C})$, there exist unique matrices $B, C \in M_n(\mathbb{C})$ such that $A = B + C$, $BC = CB$, B is diagonalizable, and C is nilpotent.*

8.16 Summary

1. *Definitions.* Given a vector space V and a linear map $T : V \rightarrow V$, T is *nilpotent* iff $T^k = 0$ for some $k \in \mathbb{N}^+$; the least such k is the *index of nilpotence* of T . T is *diagonalizable* iff for some ordered basis X of V , $[T]_X$ is a diagonal matrix. $\text{Spec}(T)$ is the set of eigenvalues of T . The *Jordan block* $J(c; k)$ is a $k \times k$ matrix with c 's on the diagonal, 1's on the next higher diagonal, and zeroes elsewhere. A *Jordan canonical form* is a matrix of the form $\text{blk-diag}(J(c_1; k_1), \dots, J(c_s; k_s))$. A *partition of n* is a non-increasing sequence $\mu = (\mu_1, \dots, \mu_s)$ of positive integers with sum n . We write $J(c; \mu) = \text{blk-diag}(J(c; \mu_1), \dots, J(c; \mu_s))$. The *diagram of μ* consists of s rows of boxes, with μ_i boxes in row i ; μ'_k is the height of the k 'th column in this diagram. Formally, $D(\mu) = \{(i, j) : 1 \leq i \leq s, 1 \leq j \leq \mu_i\}$.
2. *Nilpotent Maps and Partition Diagrams.* Given a partition μ and a basis $X = \{x(i, j) : (i, j) \in D(\mu)\}$ for V , we obtain a nilpotent linear map $T_{X, \mu}$ on V by sending $x(i, 1)$ to zero for all i , sending $x(i, j)$ to $x(i, j-1)$ for all i and

all $j > 1$, and extending by linearity. A basis for $\text{Null}(T_{X,\mu}^k)$ consists of $x(i,j)$ with $(i,j) \in D(\mu)$ and $j \leq k$; so $\dim(\text{Null}(T_{X,\mu}^k)) = \mu'_1 + \cdots + \mu'_k$. A basis for $\text{img}(T_{X,\mu}^k)$ consists of all $x(i,j)$ with $(i,j) \in D(\mu)$ and $(i,j+k) \in D(\mu)$.

3. *Classification of Nilpotent Maps.* Given any field F , any n -dimensional F -vector space V , and any nilpotent linear map $T : V \rightarrow V$, there exists a unique partition μ of n such that $T = T_{X,\mu}$ (equivalently, $[T]_X = J(0;\mu)$) for some ordered basis X of V .
4. *Jordan Canonical Form Theorem.* For all algebraically closed fields F (such as \mathbb{C}), all finite-dimensional F -vector spaces V , and all linear maps $T : V \rightarrow V$, there exists an ordered basis X of V such that $[T]_X$ is a Jordan canonical form $\text{blk-diag}(J(c_1;m_1), \dots, J(c_s;m_s))$. For any other ordered basis Y of V , if $[T]_Y = \text{blk-diag}(J(d_1;n_1), \dots, J(d_t;n_t))$, then $s = t$ and the Jordan blocks $J(d_i;n_i)$ are a rearrangement of the Jordan blocks $J(c_i;n_i)$. Every matrix $A \in M_n(F)$ is similar to a Jordan canonical form \mathbf{J} , and any other Jordan form similar to A is obtained by reordering the Jordan blocks of \mathbf{J} .
5. *Computing Jordan Forms.* The diagonal entries in any Jordan form of a linear map T (or matrix A) are the eigenvalues of T (or A). For each eigenvalue c , let $U_c = T - c\text{Id}_V$ (or $U_c = A - cI_n$). We can find an ordered basis X_c for the subspace $Z_c = \text{Null}(U_c^n)$ such that $[U_c|Z_c]_{X_c} = J(0;\mu(c))$. Letting X be the concatenation of the bases X_c , $[T]_X$ will be a Jordan canonical form. The k 'th column of $D(\mu(c))$ has size $\dim(\text{Null}(U_c^k)) - \dim(\text{Null}(U_c^{k-1}))$.
6. *Application to Differential Equations.* One way to solve $x' = Ax$ is to find an invertible P and a Jordan form \mathbf{J} with $\mathbf{J} = P^{-1}AP$. The substitution $x = Py$, $y = P^{-1}x$ converts the system $x' = Ax$ to $y' = \mathbf{J}y$, which can be solved by backsolving each Jordan block.
7. *Minimal Polynomials and Jordan Forms.* A matrix $A \in M_n(F)$ is diagonalizable iff the minimal polynomial m_A splits into distinct linear factors in $F[x]$; similarly for a linear map on an n -dimensional vector space. Given $A \in M_n(\mathbb{C})$, $m_A = \prod_{c \in \text{Spec}(A)} (x - c)^{k(c)}$ where $k(c)$ is the maximum size of any Jordan block $J(c;k)$ appearing in a Jordan form similar to A . The characteristic polynomial $\chi_A = \prod_{c \in \text{Spec}(A)} (x - c)^{m(c)}$ where $m(c)$ is the total size of all Jordan blocks $J(c;k)$ in A 's Jordan form. So, m_A divides χ_A (Cayley–Hamilton theorem). If $T \in L(V)$ is diagonalizable and W is a T -invariant subspace of V , the restriction $T|W$ is also diagonalizable.
8. *Jordan–Chevalley Decomposition.* Given a linear map T on a finite-dimensional vector space V over an algebraically closed field F , there exist unique linear maps T_d and T_n on V such that T_d is diagonalizable, T_n is nilpotent, $T_d \circ T_n = T_n \circ T_d$, and $T = T_d + T_n$.

8.17 Exercises

Unless otherwise specified, assume in these exercises that F is a field, V is a finite-dimensional F -vector space, and $T : V \rightarrow V$ is a linear map.

1. Decide whether each linear map is nilpotent. If not, explain why not. If so, find the index of nilpotence. (a) $T : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ given by $T(a,b,c) = (c,0,b)$ for $a,b,c \in \mathbb{R}$.

- (b) $T : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ given by $T(a, b, c) = (0, c, b)$ for $a, b, c \in \mathbb{R}$. (c) $D : V \rightarrow V$ given by $D(f) = df/dx$, where V is the vector space of real polynomials of degree 4 or less. (d) $T : M_3(\mathbb{R}) \rightarrow M_3(\mathbb{R})$ given by $T(A) = A^T$ (the transpose map). (e) $T : \mathbb{C}^2 \rightarrow \mathbb{C}^2$ given by $T(a, b) = (-3a + 9b, -a + 3b)$ for $a, b \in \mathbb{C}$. (f) $T : \mathbb{C}^2 \rightarrow \mathbb{C}^2$ given by $T(a, b) = (-a + 3b, -a + 3b)$ for $a, b \in \mathbb{C}$. (g) $T : \mathbb{R}^{10} \rightarrow \mathbb{R}^{10}$ given by $T(x) = Ax$ for $x \in \mathbb{R}^{10}$, where $A = \text{blk-diag}(J(0; 4), J(0; 3)^T, J(0; 3))$.
2. Decide whether each linear map is diagonalizable. If not, explain why not. If so, find an ordered basis X such that $[T]_X$ is diagonal.
 - (a) $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ given by $T(a, b) = (a, a+b)$ for $a, b \in \mathbb{R}$.
 - (b) $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ given by $T(a, b) = (a+b, a+b)$ for $a, b \in \mathbb{R}$.
 - (c) $T_\theta : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ given by $T_\theta(a, b) = (a \cos \theta - b \sin \theta, a \sin \theta + b \cos \theta)$ for $a, b \in \mathbb{R}$, where $0 < \theta < \pi/2$ is fixed.
 - (d) $T_\theta : \mathbb{C}^2 \rightarrow \mathbb{C}^2$ given by the same formula in (c) for $a, b \in \mathbb{C}$.
 - (e) $D : V \rightarrow V$ given by $D(f) = df/dx$, where V is the vector space of real polynomials of degree 4 or less.
 - (f) $T : M_2(\mathbb{R}) \rightarrow M_2(\mathbb{R})$ given by $T(A) = A^T$ (the transpose map).
 3. Fix $n > 1$.
 - (a) Prove or disprove: for all $A \in M_n(F)$, if A is nilpotent, then A is not invertible.
 - (b) Prove or disprove: for all $A \in M_n(F)$, if A is not invertible, then A is nilpotent.
 4. Use the uniqueness part of the Jordan canonical form theorem to prove that $A \in M_n(\mathbb{C})$ is diagonalizable iff every Jordan block in any Jordan canonical form similar to A has size 1.
 5.
 - (a) Suppose that $A \in M_n(F)$ has only one eigenvalue c and $A \neq cI_n$. Prove that A is not diagonalizable.
 - (b) Use (a) to explain why nonzero nilpotent maps are not diagonalizable.
 6.
 - (a) How many Jordan canonical forms are similar to the matrix $J(0; (5, 5, 5, 5, 2, 2, 2, 2, 2, 2, 2, 1, 1, 1))$?
 - (b) How many Jordan canonical forms are similar to the matrix $\text{blk-diag}(J(2; (2, 2, 2)), J(3; 3), J(4; (3, 3, 2, 2)))$?
 7. Describe the Jordan canonical form matrices that are not similar to any Jordan canonical form besides themselves.
 8.
 - (a) Suppose V is n -dimensional and $T : V \rightarrow V$ is a linear map such that for all $v \in V$, there exists $k(v) \in \mathbb{N}^+$ (depending on v) with $T^{k(v)}(v) = 0$. Prove that T is nilpotent.
 - (b) Give an example to show that the result of (a) can be false if V is infinite-dimensional.
 9. In the example at the end of §8.1, verify that Z is an ordered basis of V and $[T]_Z = J(0; (2, 1, 1))$.
 10. For each linear map T defined on an ordered basis $X = (x_1, \dots, x_n)$, decide if T is nilpotent. If it is, find an ordered basis Z and a partition μ such that $[T]_Z = J(0; \mu)$.
 - (a) $n = 5$, $T(x_i) = x_{i+1}$ for $1 \leq i < 4$, $T(x_5) = 0$.
 - (b) $n = 4$, $T(x_1) = x_3$, $T(x_2) = 0$, $T(x_3) = x_4$, $T(x_4) = x_1$.
 - (c) $n = 8$, $T(x_i) = x_{\lfloor i/2 \rfloor}$ for all i , where $x_0 = 0$.
 - (d) $n = 11$, $T(x_i) = x_{(2i \bmod 12)}$ for $1 \leq i < 12$, where $x_0 = 0$.
 11. For each linear map T defined on an ordered basis $X = (x_1, \dots, x_n)$, decide if T is nilpotent. If it is, find an ordered basis Z and a partition μ such that $[T]_Z = J(0; \mu)$.
 - (a) $n = 7$, $T(x_1) = T(x_5) = x_2$, $T(x_7) = x_4$, $T(x_2) = T(x_4) = T(x_6) = x_3$, $T(x_3) = 0$.
 - (b) $n = 3$, $T(x_1) = x_2 + 2x_3$, $T(x_2) = x_1 + 2x_3$, $T(x_3) = x_1 + 2x_2$.
 - (c) $n = 4$, $T(x_i) = x_1 + x_2 + x_3 - 3x_4$ for $1 \leq i \leq 4$.

- (d) $T(x_i) = x_{i+1} + x_{i+2} + \cdots + x_n$ for $1 \leq i < n$, $T(x_n) = 0$. (e) $n = 15$, $T(x_i) = x_{(2i \bmod 16)}$ for $1 \leq i < 15$, where $x_0 = 0$.
12. (a) Suppose $X = (x_1, \dots, x_n)$ is an ordered basis for V and T is a linear map on V such that for $1 \leq j \leq n$, either $T(x_j) = 0$ or there exists $i < j$ with $T(x_j) = x_i$. Prove T is nilpotent. (b) More generally, assume that for all j , $T(x_j) = \sum_{i=1}^{j-1} a_{ij}x_i$ for some $a_{ij} \in F$. Without using matrices or Jordan forms, prove T is nilpotent.
13. List all partitions of 5. For each partition μ , state the elements of the set $D(\mu)$, draw a picture of $D(\mu)$, and compute μ'_k for all $k \geq 1$.
14. For a certain partition μ , it is known that
- $$(\mu'_1 + \mu'_2 + \cdots + \mu'_k : k \geq 1) = (8, 12, 16, 19, 21, 23, 25, 26, 27, 28, 29, 30, 30, 30, \dots).$$
- What is μ ?
15. For a certain linear map $T_{X,\mu}$, it is known that the dimensions of $\text{img}(T_{X,\mu}^k)$ for $k \geq 1$ are: 15, 11, 8, 5, 3, 2, 1, 0. What can be said about μ ?
16. Prove: for all partitions μ of n , μ'_k is the number of $\mu_i \geq k$. Conclude that $\mu' = (\mu'_1, \mu'_2, \dots : \mu'_k > 0)$ is a partition of n , and $\mu'' = \mu$.
17. Let $X = (x_1, x_2, x_3, x_4)$ be an ordered basis of V . For each partition μ of 4, compute $T_{X,\mu}(x_i)$ for $1 \leq i \leq 4$ and $T_{X,\mu}(v)$, where $v = \sum_{i=1}^4 c_i x_i$.
18. Let $T = T_{X,\mu}$ where $\mu = (6, 3, 3, 3, 2, 1, 1)$. For all $k \geq 1$, describe a basis for $\text{Null}(T^k)$ and a basis for $\text{img}(T^k)$.
19. Let $T = T_{X,\mu}$ where $\mu = (11, 8, 6, 6, 3, 2, 2, 2, 1, 1)$. For all $k \geq 1$, compute the dimensions of $\text{Null}(T^k)$ and $\text{img}(T^k)$.
20. Assume $\dim(V) = 1$. Check directly that for any nilpotent linear map $T : V \rightarrow V$, $T = T_{(x),(1)}$ for all nonzero $x \in V$.
21. (a) Given a partition μ and $r, k \in \mathbb{N}^+$, use a visual analysis of $D(\mu)$ to find a basis for $\text{Null}(T_{X,\mu}^r) \cap \text{img}(T_{X,\mu}^k)$. Illustrate your answer for $\mu = (5, 5, 3, 1, 1, 1)$, $r = 3$, $k = 2$. (b) Similarly, find a basis for $\text{Null}(T_{X,\mu}^r) + \text{img}(T_{X,\mu}^k)$. Illustrate your answer for $\mu = (5, 5, 3, 1, 1, 1)$, $r = 1$, $k = 3$.
22. (a) Let $\mu = (5, 5, 3, 1, 1, 1)$ and $X = (x_1, \dots, x_{16})$. For each $k \geq 2$, find a partition ν and an ordered basis Z such that $T_{X,\mu}^k = T_{Z,\nu}$. (b) For a general partition μ and $k \in \mathbb{N}^+$, describe how to find Z and ν such that $T_{X,\mu}^k = T_{Z,\nu}$.
23. Suppose $[T]_X = J(0; (2, 2, 1))$ for some ordered basis $X = (x_1, \dots, x_5)$. Let $Z = (z_1, \dots, z_5)$ be a new ordered basis given by $z_j = \sum_{i=1}^5 c_{ij}x_i$ for some $c_{ij} \in F$. Find necessary and sufficient conditions on the c_{ij} 's to ensure that $[T]_Z = J(0; (2, 2, 1))$.
24. Prove: for any two subspaces W and Z of a finite-dimensional vector space V , $\dim(W + Z) + \dim(W \cap Z) = \dim(W) + \dim(Z)$. [Hint: One approach is to apply the rank-nullity theorem to the map $f : W \times Z \rightarrow W + Z$ given by $f(w, z) = w - z$ for $w \in W$ and $z \in Z$.]
25. Let $T : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ be defined by $T(a, b, c) = (2a + b - 3c, 2a + b - 3c, 2a + b - 3c)$ for $a, b, c \in \mathbb{R}$. (a) Check that T is nilpotent. (b) Find an ordered basis X_1 for $V_1 = \text{img}(T)$ and a partition ν such that $T|V_1 = T_{X_1,\nu}$. (c) Extend X_1 to an ordered basis X_2 of $V_2 = \text{img}(T) + \text{Null}(T)$, and find a partition ρ with $T|V_2 = T_{X_2,\rho}$. (d) Extend X_2 to get an indexed basis X of V and a partition μ with $T = T_{X,\mu}$. (e) Check that $X'_2 = ((1, -2, 0), (0, 3, 1))$ is a basis of V_2 such

that $T|V_2 = T_{X'_2, \rho}$. What happens if we try to solve part (d) starting from the basis X'_2 ?

26. Let $T : \mathbb{R}^4 \rightarrow \mathbb{R}^4$ be defined by

$$T(a, b, c, d) = (-a + b + c + d, a + b - c + d, -a + b + c + d, -a - b + c - d).$$

Follow the three-stage proof in the text to find X and μ such that $T = T_{X, \mu}$.

27. Let $T : \mathbb{R}^6 \rightarrow \mathbb{R}^6$ be defined by

$$T(a, b, c, d, e, f) = (0, a+2b+d-3f, 0, 2a+2b+c+d-3f, 0, 3a+2b+2c+d+e-3f).$$

Follow the three-stage proof in the text to find X and μ such that $T = T_{X, \mu}$.

[Hint: In stage 1, you can use the results of Exercise 25.]

28. Let a linear map $T : \mathbb{R}^9 \rightarrow \mathbb{R}^9$ be defined on the standard basis by $T(e_k) = e_{\lfloor k/3 \rfloor}$ for $1 \leq k \leq 9$, with $e_0 = 0$. Use the three-stage proof in the text to find a partition μ and ordered basis X with $T = T_{X, \mu}$.
29. Let $T : V \rightarrow V$ be a nilpotent linear map. Let $n = \dim(V)$ and $j = \dim(\text{Null}(T))$, so $\dim(\text{img}(T)) = n - j$. What are the possible values of $k = \dim(\text{Null}(T) \cap \text{img}(T))$? For each feasible k , construct an explicit example of a map T achieving this k .
30. *Direct Sums.* Let W_1, \dots, W_k be subspaces of V , and let X_i be an ordered basis of W_i for $1 \leq i \leq k$. Let $W = W_1 + \dots + W_k$, and let X be the concatenation of the lists X_1, \dots, X_k . Prove the following conditions are equivalent: (a) for each $w \in W$, there exist unique $w_i \in W_i$ with $w = w_1 + \dots + w_k$; (b) for $2 \leq i \leq k$, $W_i \cap (W_1 + W_2 + \dots + W_{i-1}) = \{0\}$; (c) X is an ordered basis for W ; (d) X is a linearly independent list. When any of these conditions holds, we say W is the *direct sum* of the W_i , denoted $W = W_1 \oplus W_2 \oplus \dots \oplus W_k$.
31. Give an example of three subspaces W_1, W_2, W_3 of a vector space V such that $W_1 \cap W_2 = W_1 \cap W_3 = W_2 \cap W_3 = \{0\}$, but the sum $W_1 + W_2 + W_3$ is not direct.
32. Assume $n = \dim(V) < \infty$, $S : V \rightarrow V$ is linear, and $\text{img}(S) = \text{img}(S^2)$. Must $V = \text{Null}(S) \oplus \text{img}(S)$ follow? Explain.
33. For each linear map U on $V = \mathbb{R}^4$, find explicit U -invariant subspaces Z and W satisfying the conclusions of Fitting's lemma. (a) $U(a, b, c, d) = (a, b, 0, 0)$. (b) $U(a, b, c, d) = (0, c, b, a)$. (c) $U(a, b, c, d) = (a + b, c + d, c + d, a + b)$. (d) $U(a, b, c, d) = (c - a, d - a, c - a, 2d - b - a)$.
34. Let V be the infinite-dimensional real vector space of all sequences $(x_i : i \in \mathbb{N})$ under componentwise operations. For each linear map $U : V \rightarrow V$, prove that the conclusion of Fitting's lemma does not hold for U . (a) $U((x_0, x_1, x_2, \dots)) = (x_1, x_2, x_3, \dots)$ for all $x_i \in \mathbb{R}$. (b) $U((x_0, x_1, x_2, \dots)) = (0, x_0, x_1, x_2, \dots)$ for all $x_i \in \mathbb{R}$.
35. *Eigenspaces.* For each eigenvalue c of a linear map $T : V \rightarrow V$, define the *eigenspace* $E_c = \{v \in V : T(v) = cv\}$. Thus, E_c consists of zero and all eigenvectors of T associated with the eigenvalue c . (a) Show each E_c is a subspace of V . (b) For a complex vector space V , let $X = (x_1, \dots, x_n)$ be an ordered basis of V such that $[T]_X = \text{blk-diag}(J(c; \mu), B)$, where B is a triangular matrix with no diagonal entries equal to c . Show that $(x_1, x_{\mu_1+1}, x_{\mu_1+\mu_2+1}, \dots)$ is an ordered basis of E_c . (c) Let $\text{Spec}(T) = \{c_1, \dots, c_k\}$. Use (b) to show that the sum $E_{c_1} + \dots + E_{c_k}$ is a direct sum (see Exercise 30) in the complex case. (d) For any field F , show that the sum of the eigenspaces of T is a direct sum. [Hint: If not, study a relation

- $v_1 + \cdots + v_i = 0$ with $v_j \in E_{c_j}$, $v_i \neq 0$, and i minimal.] (e) Give an example to show that the sum of all eigenspaces of V may be a proper subspace of V . If $|\text{Spec}(T)| = k$, what is the minimum possible dimension of this subspace?
36. *Generalized Eigenspaces.* Given a linear map T on an n -dimensional vector space V and $c \in \text{Spec}(T)$, the *generalized eigenspace* associated with c is $G_c = \text{Null}((T - c\text{Id}_V)^n)$. Prove: if F is algebraically closed, then $V = \bigoplus_{c \in \text{Spec}(T)} G_c$.
37. Prove: for all $A \in M_n(\mathbb{C})$, A and A^T are similar to the same Jordan canonical forms and hence are similar to one another.
38. For $c \neq 0$ and $k \geq 1$, compute $J(c; k)^{-1}$. What is a Jordan canonical form similar to this matrix?
39. Compute a Jordan canonical form similar to each matrix below.
- (a) $\begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{bmatrix}$ (b) $\begin{bmatrix} 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 \end{bmatrix}$ (c) $\begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}$
- (d) $\begin{bmatrix} 4 & -2 & 9 & -2 \\ 1 & 1 & 4 & -1 \\ 0 & 0 & 2 & 0 \\ 1 & -1 & 5 & 1 \end{bmatrix}$
40. For each matrix A in Exercise 39, find an invertible matrix P such that $P^{-1}AP$ is a Jordan canonical form.
41. Write a program that computes the Jordan canonical form of a given complex matrix.
42. Solve the following systems of ODE's. (a) $x'_1 = 3x_1$, $x'_2 = x_1 + 3x_2$; (b) $x'_1 = x_2$, $x'_2 = -x_1$; (c) $x'_1 = x_1 + 2x_2 + 3x_3$, $x'_2 = x_2 + 2x_3$, $x'_3 = x_3$; (d) $x'_1 = x_1 + x_2$, $x'_2 = x_2 + x_3$, $x'_3 = x_3 + x_4$, $x'_4 = x_1 + x_4$; (e) $x' = Ax$, where $A = \begin{bmatrix} 1 & 2 & 2 & 1 \\ 2 & 1 & 1 & 2 \\ 1 & 2 & 2 & 1 \\ 2 & 1 & 1 & 2 \end{bmatrix}$.
43. Check by induction on k that the functions in (8.12) satisfy $x' = J(c; n)x$.
44. Find the minimal polynomial and characteristic polynomial of each matrix in Exercise 39.
45. (a) Prove or disprove: for all $A \in M_n(\mathbb{C})$, if χ_A splits into distinct linear factors in $\mathbb{C}[x]$ then A is diagonalizable. (b) Prove or disprove: for all $A \in M_n(\mathbb{C})$, if A is diagonalizable then χ_A splits into distinct linear factors in $\mathbb{C}[x]$.
46. Suppose W is a T -invariant subspace of V and $h \in F[x]$ is any polynomial. Prove W is also $h(T)$ -invariant.
47. *Primary Decomposition Theorem.* For any field F and linear map $T : V \rightarrow V$, suppose $m_T = p_1^{e_1} \cdots p_s^{e_s} \in F[x]$, where p_1, \dots, p_s are distinct monic irreducible polynomials in $F[x]$. Define $W_i = \ker(p_i^{e_i}(T))$ for $1 \leq i \leq s$. (a) Show that the W_i are T -invariant subspaces of V . (b) Show that $V = W_1 + W_2 + \cdots + W_s$. [Define $g_i = m_T/p_i^{e_i} \in F[x]$ and explain why there exist $h_i \in F[x]$ with $\sum_{i=1}^s h_i g_i = 1$. For $v \in V$, show $v = \sum_{i=1}^s h_i(T)g_i(T)(v)$ where the i 'th summand is in W_i .] (c) Show that $V = W_1 \oplus W_2 \oplus \cdots \oplus W_s$ (a direct sum). [If the sum is not direct, choose t minimal such that $0 = w_1 + w_2 + \cdots + w_t$ with $w_i \in W_i$ and $w_t \neq 0$. Contradict minimality of t .]

48. Prove: for any field F , if m_T splits into a product of distinct linear factors in $F[x]$, then T is diagonalizable. (Use Exercise 47.)
49. Let S and T be diagonalizable linear maps on an F -vector space V . Prove there exists an ordered basis X such that $[S]_X$ and $[T]_X$ are both diagonal iff $S \circ T = T \circ S$. (So two diagonalizable operators can be *simultaneously diagonalized* iff the operators commute.) [Hint: For the backward direction, choose an ordered basis Y such that $[S]_Y = \text{blk-diag}(c_1 I_{n_1}, \dots, c_s I_{n_s})$, where c_1, \dots, c_s are the distinct eigenvalues of S . Break Y into sublists Y_1, \dots, Y_s of length n_1, \dots, n_s , and let W_1, \dots, W_s be the subspaces spanned by these sublists. Explain why each W_i is T -invariant. Show that there is an ordered basis X_i for W_i such that $[S|W_i]_{X_i}$ and $[T|W_i]_{X_i}$ are both diagonal.]
50. (a) Prove: if S and T are diagonalizable linear maps on V and $S \circ T = T \circ S$, then $aS + bT$ is diagonalizable for any $a, b \in F$. (b) Give a specific example of diagonalizable linear maps S and T such that $S - T$ is not diagonalizable.
51. (a) Suppose S and T are commuting nilpotent linear maps on an F -vector space V . Prove $S - T$ is nilpotent. (b) Must (a) hold if S and T do not commute? Prove or give a counterexample.
52. Deduce the Jordan–Chevalley decomposition for matrices from the corresponding result for linear maps.
53. Find the Jordan–Chevalley decomposition of each matrix in Exercise 39.
54. Find the Jordan–Chevalley decomposition of the matrix A in (8.11).
55. Given $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$, show that there exist infinitely many pairs of matrices B, C with $A = B + C$, B diagonalizable, and C nilpotent. Why doesn't this contradict the uniqueness assertion in the Jordan–Chevalley theorem?
56. Given $T \in L(V)$, show that there exist polynomials $f, g \in F[x]$ with no constant term such that $T_d = f(T)$ and $T_n = g(T)$.

Matrix Factorizations

Matrices with special structure — such as diagonal matrices, triangular matrices, and unitary matrices — are simpler to work with than general matrices. Many algorithms have been developed in numerical linear algebra to convert an input matrix into a form with specified special structure by using a sequence of carefully chosen matrix operations. These algorithms can often be described mathematically as providing *factorizations* of matrices into products of structured matrices. This chapter proves the existence and uniqueness properties of several matrix factorizations and explores the algebraic and geometric ideas leading to these factorizations.

The first factorization, called the QR factorization, writes a complex matrix as a product of a unitary matrix Q and an upper-triangular matrix R . More generally, Q can be replaced by a rectangular matrix with orthonormal columns. To obtain such factorizations, we study the Gram–Schmidt algorithm for converting a linearly independent list of vectors into an orthonormal list of vectors. Another approach to QR factorizations involves Householder transformations, which generalize geometric reflections in \mathbb{R}^2 or \mathbb{R}^3 . By cleverly choosing and composing Householder matrices, we can force the entries of a matrix below the main diagonal to become zero, one column at a time.

The second family of factorizations, called LU decompositions, writes certain square matrices as products LU with L lower-triangular and U upper-triangular. When such a product exists, we can arrange for one of the matrices (L or U) to have all ones on its main diagonal. In general, a triangular matrix is called *unitriangular* iff all of its diagonal entries are equal to one. Not every matrix has an LU factorization, but we will see that an appropriate permutation of the rows will convert any invertible matrix to a matrix of the form LU . Similarly, by permuting both rows and columns, any matrix can be converted to a matrix with an LU factorization. These LU factorizations are closely connected to the Gaussian elimination algorithm, which reduces a matrix to a simpler form by systematically creating zero entries.

One motivation for finding QR and LU factorizations is the efficient solution of a linear system of equations $Ax = b$. If A is lower-triangular, this system is readily solved by *forward substitution*, in which we solve for the components x_1, \dots, x_n of x in this order. Similarly, if A is upper-triangular, the system can be solved quickly by *backward substitution*, in which we solve for x_n, \dots, x_1 in this order. If $A = LU$ with L lower-triangular and U upper-triangular, we can solve $Ax = LUX = b$ by first solving $Ly = b$ for y via forward substitution, and then solving $Ux = y$ for x via backward substitution. Given $A = QR$ with Q unitary and R upper-triangular, we can solve $Ax = QRx = b$ by noting that $Q^{-1} = Q^*$ (which can be computed quickly from Q), so $Ax = b$ iff $Rx = Q^*b$. The latter system can be solved by backward substitution.

The chapter concludes with two more matrix factorization results. The Cholesky factorization expresses a positive semidefinite matrix as a product LL^* , where L is lower-triangular. This factorization is closely related to the QR and LU factorizations, and it has applications to least-squares approximation problems. The singular value decomposition expresses a matrix A as a product of a unitary matrix, a diagonal matrix with nonnegative

diagonal entries, and another unitary matrix. These factorizations all play a fundamental role in numerical linear algebra.

9.1 Approximation by Orthonormal Vectors

To prepare for our study of the QR factorization, we first discuss how to approximate a vector as a linear combination of orthonormal vectors. Let V be a complex inner product space. Recall that this is a complex vector space with an *inner product* (the analogue of the dot product in \mathbb{R}^n) satisfying these conditions for all $x, y, z \in V$ and $c \in \mathbb{C}$: $\langle x + y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$; $\langle cx, y \rangle = c\langle x, y \rangle$; $\langle y, x \rangle = \overline{\langle x, y \rangle}$; and for $x \neq 0$, $\langle x, x \rangle \in \mathbb{R}^+$. The most frequently used complex inner product space is \mathbb{C}^n with the “standard” inner product $\langle v, w \rangle = \sum_{k=1}^n v_k \overline{w_k}$.

For any $x \in V$, the *length* of x is $\|x\| = \sqrt{\langle x, x \rangle}$. The *distance* between vectors $x, y \in V$ is $\|x - y\|$. Two vectors $x, y \in V$ are called *orthogonal* iff $\langle x, y \rangle = 0$ iff $\langle y, x \rangle = 0$. Orthogonal vectors x and y satisfy the *Pythagorean theorem* $\|x + y\|^2 = \|x\|^2 + \|y\|^2$, since

$$\|x + y\|^2 = \langle x + y, x + y \rangle = \langle x, x \rangle + \langle x, y \rangle + \langle y, x \rangle + \langle y, y \rangle = \|x\|^2 + 0 + 0 + \|y\|^2.$$

A list of vectors (u_1, \dots, u_k) in V is called *orthogonal* iff $\langle u_r, u_s \rangle = 0$ for all $r \neq s$ in $[k]$. The list is called *orthonormal* iff the list is orthogonal and $\langle u_r, u_r \rangle = 1$ for all $r \in [k]$. Geometrically, an orthonormal list consists of mutually perpendicular unit vectors in V .

An *orthogonal list of nonzero vectors* (u_1, \dots, u_k) must be linearly independent. For suppose $c_1 u_1 + \dots + c_k u_k = 0$ for some $c_1, \dots, c_k \in \mathbb{C}$. Fix $r \in [k]$ and take the inner product of this linear combination with u_r . We get

$$0 = \langle 0, u_r \rangle = \langle c_1 u_1 + \dots + c_k u_k, u_r \rangle = \sum_{j=1}^k c_j \langle u_j, u_r \rangle = c_r \langle u_r, u_r \rangle.$$

Since u_r is nonzero, $\langle u_r, u_r \rangle > 0$, and so $c_r = 0$. In particular, an orthonormal list (u_1, \dots, u_k) is linearly independent.

Now consider the following *approximation problem*: given an orthonormal list (u_1, \dots, u_k) in the complex inner product space V and any vector $v \in V$, find complex scalars c_1, \dots, c_k such that $x = c_1 u_1 + \dots + c_k u_k$ has minimum distance to v among all vectors in the subspace U spanned by (u_1, \dots, u_k) . We will show there exists a unique $x \in U$ minimizing $\|v - x\|$, which is obtained by choosing $c_r = \langle v, u_r \rangle$ for $r \in [k]$, and moreover $v - x$ is orthogonal to every vector in U . We call x the *orthogonal projection* of v onto U .

Define $x = \sum_{r=1}^k \langle v, u_r \rangle u_r \in U$. Given any $s \in [k]$, we compute

$$\langle v - x, u_s \rangle = \langle v, u_s \rangle - \langle x, u_s \rangle = \langle v, u_s \rangle - \sum_{r=1}^k \langle v, u_r \rangle \langle u_r, u_s \rangle = \langle v, u_s \rangle - \langle v, u_s \rangle = 0.$$

Given any $u \in U$, write $u = \sum_{s=1}^k d_s u_s$ with $d_s \in \mathbb{C}$, and note

$$\langle v - x, u \rangle = \sum_{s=1}^k \overline{d_s} \langle v - x, u_s \rangle = \sum_{s=1}^k \overline{d_s} \cdot 0 = 0.$$

Hence, $v - x$ is orthogonal to every $u \in U$. In particular, for any $y \in U$, $x - y$ lies in the subspace U , so the Pythagorean theorem gives

$$\|v - y\|^2 = \|(v - x) + (x - y)\|^2 = \|v - x\|^2 + \|x - y\|^2.$$

Since $\|x-y\|^2 = \langle x-y, x-y \rangle \geq 0$, with equality iff $x-y=0$, we see that $\|v-y\|^2 \geq \|v-x\|^2$ for all $y \in U$, with equality iff $y=x$. So x as defined above is the unique vector in U with minimum distance to v .

The above discussion shows that any $v \in V$ can be decomposed into a sum $v = x + (v-x)$, where $x = \sum_{r=1}^k \langle v, u_r \rangle u_r \in U$ and $v-x$ is orthogonal to every vector in U . Observe that v lies in the subspace U spanned by (u_1, \dots, u_r) iff v itself is the unique vector in U with minimum distance to v iff $x=v$ iff $v-x=0$ iff $v = \sum_{r=1}^k \langle v, u_r \rangle u_r$. This gives an algorithm for detecting when a given $v \in V$ is a linear combination of the orthonormal list (u_1, \dots, u_k) and finding the coefficients of the linear combination when it exists.

For example, consider $u_1 = 0.5(1, 1, 1, 1)$, $u_2 = 0.5(1, -1, 1, -1)$, and $u_3 = 0.5(-1, -1, 1, 1)$ in $V = \mathbb{C}^4$. One readily checks that (u_1, u_2, u_3) is an orthonormal list. Given $v = (-5, 1, 3, 9)$, we compute $\langle v, u_1 \rangle = 4$, $\langle v, u_2 \rangle = -6$, and $\langle v, u_3 \rangle = 8$. Then $x = 4u_1 - 6u_2 + 8u_3 = v$, so $v \in U$ and we have written v as a specific linear combination of u_1, u_2, u_3 . For $v = (2, -3, 0, 1)$, we compute $\langle v, u_1 \rangle = 0$, $\langle v, u_2 \rangle = 2$, and $\langle v, u_3 \rangle = 1$. Here $x = 0u_1 + 2u_2 + 1u_3 = (0.5, -1.5, 1.5, -0.5)$ is the orthogonal projection of v onto U . We can see directly that $v-x = (1.5, -1.5, -1.5, 1.5)$ is orthogonal to x, u_1, u_2, u_3 , and hence to every vector in U . Although $v-x$ is not a unit vector, we can change it into one by dividing by its length, obtaining $u_4 = 0.5(1, -1, -1, 1)$. We now have a longer orthonormal list (u_1, u_2, u_3, u_4) , which (being linearly independent) must be a basis of \mathbb{C}^4 . This illustrates the basic step in the Gram–Schmidt orthonormalization algorithm, which we describe in the next section.

9.2 Gram–Schmidt Orthonormalization

Given a complex inner product space V , a finite-dimensional subspace U , and a vector $v \in V$, our test for determining whether v is in the subspace U requires us to know an orthonormal basis of U . This suggests the question of how to convert an *arbitrary* basis for a subspace U into an *orthonormal* basis of U .

Let (v_1, \dots, v_k) be any linearly independent list of vectors in V . We describe a process, called the *Gram–Schmidt orthonormalization algorithm*, that will produce an orthonormal list (u_1, \dots, u_k) such that for $1 \leq r \leq k$, the sublist (v_1, \dots, v_r) spans the same subspace as the sublist (u_1, \dots, u_r) . The classical version of the algorithm executes the following steps for $s = 1, 2, \dots, k$ in turn: at step s , u_1, \dots, u_{s-1} have already been found; calculate $x_s = \sum_{r=1}^{s-1} \langle v_s, u_r \rangle u_r$ (taking $x_1 = 0$); and set $u_s = (v_s - x_s)/\|v_s - x_s\|$. Intuitively, we obtain the next vector u_s by subtracting from v_s its orthogonal projection x_s onto the span of the preceding vectors, and then normalizing to get a unit vector.

To see that the algorithm has the required properties, we prove by induction on $r \in [k]$ that (u_1, \dots, u_r) is an orthonormal basis for the subspace V_r spanned by (v_1, \dots, v_r) . For the base case $r = 1$, note that $x_1 = 0$, $v_1 \neq 0$, and so $u_1 = v_1/\|v_1\|$ is a well-defined vector of length 1. Because u_1 is a nonzero scalar multiple of v_1 , (u_1) is a basis of V_1 .

For the induction step, fix r with $1 < r \leq k$, and assume we already know (u_1, \dots, u_{r-1}) is an orthonormal basis of V_{r-1} . On one hand, since (v_1, \dots, v_r) is linearly independent by hypothesis, v_r is not in the subspace V_{r-1} spanned by (v_1, \dots, v_{r-1}) . Since (u_1, \dots, u_{r-1}) spans this subspace by assumption, the orthogonal projection x_r is not equal to v_r . Thus $v_r - x_r \neq 0$, and it makes sense to divide this vector by its length to get a unit vector u_r . On the other hand, we proved earlier that $v_r - x_r$ is orthogonal to every vector in the space V_{r-1} , and the same is true of the scalar multiple u_r of $v_r - x_r$. In particular, u_r

is orthogonal to u_1, \dots, u_{r-1} , proving orthonormality (and hence linear independence) of the list (u_1, \dots, u_r) . Finally, since x_r lies in V_{r-1} and $u_r = ||v_r - x_r||^{-1}(v_r - x_r)$, we see that u_r lies in the subspace V_r spanned by (v_1, \dots, v_r) . Since the r linearly independent vectors u_1, \dots, u_r all belong to the r -dimensional subspace V_r , these vectors must be an orthonormal basis for this subspace, completing the induction.

For example, suppose we are given the linearly independent vectors

$$v_1 = (2, 2, 2, 2), \quad v_2 = (4, -1, 4, -1), \quad v_3 = (-4, -2, 0, 2), \quad v_4 = (2, -3, 0, 1)$$

in $V = \mathbb{C}^4$. Executing the Gram–Schmidt algorithm, we first find that $x_1 = 0$ and $u_1 = v_1/||v_1|| = (0.5, 0.5, 0.5, 0.5)$. Second, $x_2 = \langle v_2, u_1 \rangle u_1 = 3u_1 = (1.5, 1.5, 1.5, 1.5)$, $v_2 - x_2 = (2.5, -2.5, 2.5, -2.5)$, and $u_2 = (v_2 - x_2)/||v_2 - x_2|| = (0.5, -0.5, 0.5, -0.5)$. Third, $x_3 = \langle v_3, u_1 \rangle u_1 + \langle v_3, u_2 \rangle u_2 = -2u_1 - 2u_2 = (-2, 0, -2, 0)$, $v_3 - x_3 = (-2, -2, 2, 2)$, and $u_3 = (v_3 - x_3)/||v_3 - x_3|| = (-0.5, -0.5, 0.5, 0.5)$. Finally, as in the example at the end of §9.1, we find $x_4 = 0u_1 + 2u_2 + 1u_3$, $v_4 - x_4 = (1.5, -1.5, -1.5, 1.5)$, and $u_4 = (0.5, -0.5, -0.5, 0.5)$. So (u_1, u_2, u_3, u_4) is an orthonormal basis of \mathbb{C}^4 .

We close with a few remarks about the Gram–Schmidt orthonormalization process. First, this algorithm provides a constructive proof of the theorem that *every subspace of a finite-dimensional complex inner product space has an orthonormal basis*. Second, one can use the algorithm to see that *every orthonormal list in a finite-dimensional space V can be extended to an orthonormal ordered basis of V* (Exercise 12). Third, the algorithm can also be applied to real inner product spaces by restricting to real scalars. Fourth, the formulas in the algorithm remain valid for an infinite-dimensional complex inner product space, transforming a (countably) infinite linearly independent sequence $(v_1, v_2, \dots, v_n, \dots)$ into an infinite orthonormal sequence $(u_1, u_2, \dots, u_n, \dots)$ such that (v_1, \dots, v_n) and (u_1, \dots, u_n) span the same subspace for all $n \in \mathbb{N}^+$. In this situation, however, the “algorithm” does not terminate in finitely many steps.

Finally, we can even apply the algorithm to a list of vectors (v_1, \dots, v_k) that may be linearly dependent. If the list is linearly dependent and r is minimal such that v_r is a linear combination of v_1, \dots, v_{r-1} , then the algorithm will detect this by computing that the orthogonal projection x_r is equal to v_r , hence $v_r - x_r = 0$. In this case, we leave u_r undefined, discard v_r , and continue processing v_{r+1}, v_{r+2} , and so on. If further linear dependencies exist, they will be detected as they arise. In the end, the algorithm tells us exactly which v_r 's depend on preceding v_j 's and provides an orthonormal list (of size k or less) that spans the same subspace as the list (v_1, \dots, v_k) .

9.3 Gram–Schmidt QR Factorization

We now recast our results on the Gram–Schmidt orthonormalization algorithm as a theorem about matrix factorizations. We will show that *for any matrix $A \in M_{n,k}(\mathbb{C})$ with k linearly independent columns, there exist a unique matrix $Q \in M_{n,k}(\mathbb{C})$ with orthonormal columns and a unique upper-triangular matrix $R \in M_k(\mathbb{C})$ with all diagonal entries in \mathbb{R}^+ , such that $A = QR$. When A is a real matrix, Q and R are real also.*

To prove existence of the factorization, we consider the complex inner product space \mathbb{C}^n (viewed as a set of column vectors) and look at the linearly independent list of vectors (v_1, \dots, v_k) , where $v_j = A^{[j]}$ is the j 'th column of the given matrix A . The orthonormalization algorithm produces an orthonormal list (u_1, \dots, u_k) such that, for $1 \leq j \leq k$, the subspace V_j spanned by (v_1, \dots, v_j) coincides with the subspace spanned by (u_1, \dots, u_j) . Let $Q \in M_{n,k}(\mathbb{C})$ be the matrix with columns $Q^{[j]} = u_j \in \mathbb{C}^n$ for $1 \leq j \leq k$.

Define an upper-triangular matrix $R \in M_k(\mathbb{C})$ by setting $R(i, j) = \langle v_j, u_i \rangle$ for $1 \leq i \leq j \leq k$, and $R(i, j) = 0$ for all $i > j$. Note that Q and R have real entries if A does.

To check that $A = QR$, we use facts established in §4.7 and §9.1. First, the j 'th column of QR is $Q(R^{[j]})$. Second, $Q(R^{[j]})$ is a linear combination of the columns of Q with coefficients given by the entries of $R^{[j]}$. Specifically, the j 'th column of QR is $R(1, j)Q^{[1]} + R(2, j)Q^{[2]} + \cdots + R(k, j)Q^{[k]}$. Using our definitions of Q and R , the j 'th column of QR is

$$\langle v_j, u_1 \rangle u_1 + \langle v_j, u_2 \rangle u_2 + \cdots + \langle v_j, u_j \rangle u_j.$$

Since v_j is in V_j , which is the span of (u_1, \dots, u_j) , the linear combination just written must be equal to v_j (see §9.1). This means that the j 'th column of QR is $v_j = A^{[j]}$ for all j , so that $A = QR$ as needed.

We see that the diagonal entries of R are strictly positive real numbers as follows. Recall from the description of the orthonormalization algorithm that $u_j = (v_j - x_j)/\|v_j - x_j\|$, where x_j is the orthogonal projection of v_j onto the subspace V_{j-1} . We know that $v_j - x_j$ is orthogonal to everything in V_{j-1} , hence is orthogonal to x_j . Then $\|v_j - x_j\|^2 = \langle v_j - x_j, v_j - x_j \rangle = \langle v_j, v_j - x_j \rangle - \langle x_j, v_j - x_j \rangle = \langle v_j, v_j - x_j \rangle$. So, the diagonal entry $R(j, j)$ is $\langle v_j, u_j \rangle = \langle v_j, v_j - x_j \rangle / \|v_j - x_j\| = \|v_j - x_j\|^2 / \|v_j - x_j\| = \|v_j - x_j\| \in \mathbb{R}^+$.

Before discussing uniqueness, we consider an example. Suppose $A = \begin{bmatrix} 2 & 4 & -4 \\ 2 & -1 & -2 \\ 2 & 4 & 0 \\ 2 & -1 & 2 \end{bmatrix}$.

By the calculations in §9.2 and the definitions above, $A = QR$ holds for

$$Q = \begin{bmatrix} 1/2 & 1/2 & -1/2 \\ 1/2 & -1/2 & -1/2 \\ 1/2 & 1/2 & 1/2 \\ 1/2 & -1/2 & 1/2 \end{bmatrix}, \quad R = \begin{bmatrix} 4 & 3 & -2 \\ 0 & 5 & -2 \\ 0 & 0 & 4 \end{bmatrix}.$$

Returning to the general case, we now prove uniqueness of Q and R . Assume $A = QR = Q_1 R_1$, where $Q, Q_1 \in M_{n,k}(\mathbb{C})$ both have orthonormal columns and $R, R_1 \in M_k(\mathbb{C})$ are both upper-triangular with all diagonal entries in \mathbb{R}^+ . As above, let the columns of A be v_1, \dots, v_k , let the columns of Q be u_1, \dots, u_k , and let the columns of Q_1 be z_1, \dots, z_k . We show by strong induction on $j \in [k]$ that $u_j = z_j$ and $R^{[j]} = R_1^{[j]}$. Fix $j \in [k]$, and assume $u_s = z_s$ and $R^{[s]} = R_1^{[s]}$ is already known for all s with $1 \leq s < j$. On one hand, consideration of the j 'th column of $A = QR$ shows (as in the existence proof) that

$$v_j = R(1, j)u_1 + \cdots + R(j-1, j)u_{j-1} + R(j, j)u_j.$$

Taking the inner product of each side with u_i , we see that $R(i, j) = \langle v_j, u_i \rangle$ for $1 \leq i \leq j$. Since R is upper-triangular, $R(i, j) = 0$ for $j < i \leq k$.

On the other hand, looking at the j 'th column of $A = Q_1 R_1$ gives

$$\begin{aligned} v_j &= R_1(1, j)z_1 + \cdots + R_1(j-1, j)z_{j-1} + R_1(j, j)z_j \\ &= R_1(1, j)u_1 + \cdots + R_1(j-1, j)u_{j-1} + R_1(j, j)z_j. \end{aligned}$$

Since the list $(z_1, \dots, z_{j-1}, z_j) = (u_1, \dots, u_{j-1}, z_j)$ is orthonormal, taking the inner product of each side with u_i gives $R_1(i, j) = \langle v_j, u_i \rangle = R(i, j)$ for $1 \leq i < j$. Since R_1 is upper-triangular, $R_1(i, j) = 0$ for $j < i \leq k$. Now, equating the two expressions for v_j and cancelling $\sum_{i=1}^{j-1} \langle v_j, u_i \rangle u_i$, we deduce that $R(j, j)u_j = R_1(j, j)z_j$. Taking the length of both sides gives $|R(j, j)| = |R_1(j, j)|$. Since both $R(j, j)$ and $R_1(j, j)$ are positive real numbers, we conclude finally that $R(j, j) = R_1(j, j)$ and $u_j = z_j$. This completes the inductive proof of uniqueness.

We can extend the QR factorization to arbitrary matrices as follows. *Given any matrix $A \in M_{n,k}(\mathbb{C})$ with column rank r , there exist $Q \in M_{n,r}(\mathbb{C})$ with orthonormal columns and an upper-triangular $R \in M_{r,k}(\mathbb{C})$ with all diagonal entries nonnegative real numbers such that $A = QR$.* (A possibly non-square matrix R is *upper-triangular* iff $R(i,j) = 0$ for all $i > j$.) As before, let the columns of A be v_1, \dots, v_k . Let the columns of Q be the orthonormal list (u_1, \dots, u_r) produced when the Gram–Schmidt algorithm is applied to the list (v_1, \dots, v_k) . We have $r < k$ whenever some of the v_j 's are linear combinations of previous v_i 's. Let $R(i,j) = \langle v_j, u_i \rangle$ for $1 \leq i \leq j \leq k$, and $R(i,j) = 0$ for $r \geq i > j \geq 1$. As in the earlier existence proof, $A = QR$ follows from the fact that each v_j is in the span of the orthonormal list $(u_1, \dots, u_{\min(j,r)})$. If some v_j is in the span of (u_1, \dots, u_{j-1}) , then $R(j,j)$ will be zero. By deleting the columns of A corresponding to those v_j 's that depend linearly on earlier v_i 's, and deleting the same columns of R , we obtain a factorization $A' = QR'$ of the type initially discussed, with $A' \in M_{n,r}(\mathbb{C})$ having full column rank and $R' \in M_{r,r}(\mathbb{C})$ being square and upper-triangular with all strictly positive diagonal entries.

9.4 Householder Reflections

Our next goal is to derive a version of the QR factorization in which an arbitrary matrix $A \in M_{n,k}(\mathbb{C})$ is factored as $A = QR$, where $Q \in M_n(\mathbb{C})$ is a unitary matrix and $R \in M_{n,k}(\mathbb{C})$ is an upper-triangular matrix. One can obtain such a factorization by modifying the proof given earlier based on Gram–Schmidt orthonormalization (Exercise 16). This section and the next one describe Householder's algorithm for reaching this factorization, which has certain computational advantages over algorithms using Gram–Schmidt.

The key idea is that we can apply a sequence of reflections to transform A into an upper-triangular matrix by forcing the required zeroes to appear, one column at a time. In \mathbb{R}^n , a *reflection* is a linear map that sends a given nonzero vector v to $-v$ and fixes all vectors w perpendicular to v . The subspace of all such vectors w , which is the orthogonal complement of $\{v\}$, is the “mirror” through which v is being reflected.

Let us define reflections formally in the complex inner product space \mathbb{C}^n with the standard inner product $\langle x, y \rangle = \sum_{k=1}^n x_k \bar{y}_k = y^* x$ for $x, y \in \mathbb{C}^n$. Fix a nonzero vector $v \in \mathbb{C}^n$. Starting with any basis (v, v_2, \dots, v_n) of \mathbb{C}^n , we can use Gram–Schmidt to obtain an orthonormal ordered basis $B = (u_1, u_2, \dots, u_n)$ of \mathbb{C}^n , where $u_1 = v/\|v\|$. There is a unique linear map $T_v : \mathbb{C}^n \rightarrow \mathbb{C}^n$ defined by setting

$$T_v(c_1 u_1 + c_2 u_2 + \cdots + c_n u_n) = -c_1 u_1 + c_2 u_2 + \cdots + c_n u_n \quad (9.1)$$

for all $c_1, \dots, c_n \in \mathbb{C}$. Observe that T_v sends any scalar multiple of u_1 to its negative, whereas T_v sends any linear combination of u_2, \dots, u_n to itself. In particular, $T_v(v) = -v$ and $T_v(w) = w$ for all w orthogonal to v ; these formulas uniquely determine the linear map T_v and show that T_v does not depend on the choice of the orthonormal basis B . The matrix of T_v relative to the ordered basis B is diagonal with diagonal entries $-1, 1, \dots, 1$. This matrix is evidently unitary and equal to its own inverse, so $\|T_v(z)\| = \|z\|$ for all $z \in \mathbb{C}^n$, and $T_v^{-1} = T_v$.

Define a matrix $Q_v = I_n - (2/\|v\|^2)vv^* \in M_n(\mathbb{C})$. We claim that $T_v(y) = Q_v y$ for all $y \in \mathbb{C}^n$. By linearity, it suffices to verify this for $y = v, u_2, \dots, u_n$. For $y = v$, recall that $\|v\|^2 = \langle v, v \rangle = v^* v$ and compute

$$Q_v v = I_n v - \left[\frac{2}{\|v\|^2} vv^* \right] v = v - \frac{2}{v^* v} v(v^* v) = v - 2v = -v = T_v(v).$$

For $y = u_j$ with $2 \leq j \leq n$, note $v^* u_j = \langle u_j, v \rangle = \langle u_j, ||v|| u_1 \rangle = 0$, so

$$Q_v u_j = I_n u_j - \left[\frac{2}{||v||^2} v v^* \right] u_j = u_j - \frac{2}{||v||^2} v (v^* u_j) = u_j = T_v(u_j).$$

So the claim holds. We call T_v the *Householder transformation determined by v* and Q_v the *Householder matrix determined by v* . We also refer to T_v and Q_v as *Householder reflections* to emphasize the geometric character of these maps and matrices.

One may check from the explicit formula that Q_v is unitary ($Q_v^* Q_v = Q_v Q_v^* = I_n$) and self-inverse ($Q_v^{-1} = Q_v$); these facts also follow from the corresponding properties of T_v or $[T_v]_B$. It is also routine to confirm that for any nonzero scalar $c \in \mathbb{C}$, $T_{cv} = T_v$ and $Q_{cv} = Q_v$. Note that if v is a unit vector, the formula for Q_v simplifies to $Q_v = I_n - 2vv^*$, and we can always arrange this by replacing v by an appropriate scalar multiple of itself. We also define $T_0 = \text{id}_{\mathbb{C}^n}$ and $Q_0 = I_n$.

The following lemma shows how to find Householder reflections sending a given vector x to certain other vectors having the same length as x . *For any $x, y \in \mathbb{C}^n$ with $x \neq y$ and $||x|| = ||y||$ and $\langle x, y \rangle \in \mathbb{R}$, there exists $v \in \mathbb{C}^n$ with $Q_v x = y$; specifically, $Q_v x = y$ iff v is a nonzero scalar multiple of $x - y$.* To prove the lemma, fix $x \neq y$ in \mathbb{C}^n with $||x|| = ||y||$ and $\langle x, y \rangle \in \mathbb{R}$. Define $v = x - y$ and $w = x + y$; the geometric motivation for choosing v and w in this way is suggested in Figure 9.1. The hypothesis $\langle x, y \rangle \in \mathbb{R}$ implies $\langle x, y \rangle = \langle y, x \rangle = \langle y, x \rangle$, so that

$$\langle v, w \rangle = \langle x - y, x + y \rangle = \langle x, x \rangle - \langle y, x \rangle + \langle x, y \rangle - \langle y, y \rangle = ||x||^2 - ||y||^2 = 0.$$

Since w is orthogonal to v , $T_v(w) = w$. Now $x = (x - y)/2 + (x + y)/2 = v/2 + w/2$, so

$$T_v(x) = T_v(v/2 + w/2) = \frac{1}{2}T_v(v) + \frac{1}{2}T_v(w) = (-v)/2 + w/2 = (y - x)/2 + (x + y)/2 = y$$

(cf. Figure 9.1). Thus, $Q_v x = y$, and hence $Q_{cv} x = y$ for any nonzero $c \in \mathbb{C}$.

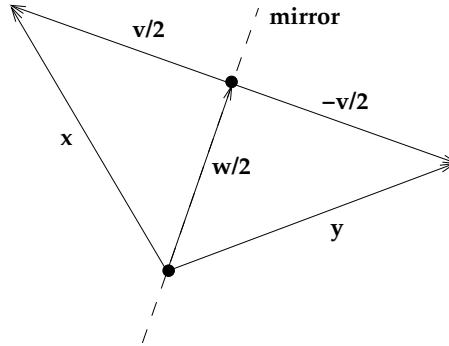


FIGURE 9.1

Finding a Reflection that Sends x to y .

To prove the converse, assume $z \in \mathbb{C}^n$ satisfies $T_z(x) = y$. Then $T_z(y) = T_z^{-1}(y) = x$, so $T_z(v) = T_z(x - y) = T_z(x) - T_z(y) = y - x = -v$. From the defining formula for T_z (see (9.1) with v there replaced by z), we find that the only vectors sent to their negatives by T_z are scalar multiples of z . So v is a nonzero scalar multiple of z and vice versa, completing the proof of the lemma. Observe that v is not unique even if we insist that it be a unit vector, since we can still replace v by $e^{i\theta}v$ for any $\theta \in [0, 2\pi)$. In order for v to exist, we must assume $||x|| = ||y||$ and $\langle x, y \rangle \in \mathbb{R}$ (see Exercise 24).

9.5 Householder QR Factorization

Given $A \in M_{n,k}(\mathbb{C})$, *Householder's QR algorithm* finds a factorization $A = QR$, where $Q \in M_n(\mathbb{C})$ is unitary and $R \in M_{n,k}(\mathbb{C})$ is upper-triangular with nonnegative diagonal entries, by determining a sequence of Householder reflections and a diagonal unitary matrix whose product is Q . The main idea is to use the lemma of §9.4 repeatedly to create the zero entries in R one column at a time.

To begin, let $x = A^{[1]} \in \mathbb{C}^n$ be the first column of A . Write $x(1) = A(1, 1) = se^{i\theta}$ for some nonnegative real s and $\theta \in [0, 2\pi)$. Define $y = -||x||e^{i\theta}e_1 \in \mathbb{C}^n$, and note that $||y|| = ||x||$ and $\langle x, y \rangle = -se^{i\theta}||x||e^{-i\theta}$ is real. (Exercise 31 discusses why the negative sign is present in the definition of y .) Let $v_1 = x - y$, and consider the matrix $Q_{v_1}A$. The first column of this matrix is $Q_{v_1}x = y$, which is a multiple of e_1 and hence has zeroes in rows 2 through n . Since $Q_0 = I_n$, this conclusion holds even if $x = 0$ or $x = y$.

At the next stage, let A_2 be the submatrix of $Q_{v_1}A$ obtained by deleting the first row and column. Repeat the construction in the previous paragraph, taking $x \in \mathbb{C}^{n-1}$ to be the first column of A_2 , to obtain $v'_2 \in \mathbb{C}^{n-1}$ such that $Q_{v'_2}A_2 \in M_{n-1,k-1}(\mathbb{C})$ has zeroes in the first column below row 1. Construct $v_2 \in \mathbb{C}^n$ by preceding v'_2 with a zero. The first row and column of $v_2v_2^*$ contain all zeroes, so Q_{v_2} is a block-diagonal matrix with diagonal blocks I_1 and $Q_{v'_2}$. It follows that $Q_{v_2}Q_{v_1}A$ has only zeroes below the diagonal in the first two columns.

We continue similarly to process the remaining columns. After j steps, we have found a matrix $Q_{v_j} \cdots Q_{v_1}A$ with zeroes below the diagonal in the first j columns. Repeat the construction in step 1, taking $x \in \mathbb{C}^{n-j}$ to be the first column of the lower-right $(n-j) \times (k-j)$ submatrix of the current matrix. We thereby obtain $v'_{j+1} \in \mathbb{C}^{n-j}$ such that $Q_{v'_{j+1}}$ sends x to a multiple of $e_1 \in \mathbb{C}^{n-j}$. Let $v_{j+1} \in \mathbb{C}^n$ be v'_{j+1} preceded by j zeroes, so that $Q_{v_{j+1}}$ is block-diagonal with diagonal blocks I_j and $Q_{v'_{j+1}}$. Then $Q_{v_{j+1}}Q_{v_j} \cdots Q_{v_1}A$ has only zeroes below the diagonal in the first $j+1$ columns.

After $s = \min(n-1, k)$ steps, we have a relation $Q_{v_s} \cdots Q_{v_1}A = R'$, where $R' \in M_{n,k}(\mathbb{C})$ is upper-triangular. Left-multiplying both sides by an appropriate diagonal matrix D with all diagonal entries of the form $e^{i\theta}$, we can obtain a factorization $DQ_{v_s} \cdots Q_{v_1}A = R$ where R is upper-triangular with nonnegative real diagonal entries. Solving for A and recalling that $Q_{v_j}^{-1} = Q_{v_j}$ for all j , we get $A = QR$ with $Q = Q_{v_1}Q_{v_2} \cdots Q_{v_s}D^{-1}$. The matrix Q is unitary, being a product of unitary matrices.

We illustrate Householder's QR algorithm on the matrix

$$A = \begin{bmatrix} -5 & 3 & -9 & 1 \\ -4 & 1 & 3 & 12 \\ 8 & 3 & -2 & 0 \\ 4 & 9 & 1 & -4 \end{bmatrix}.$$

In the first step, $x = (-5, -4, 8, 4)^T$, $||x|| = 11$, $x(1) = -5 = 5e^{i\pi}$, so $y = -||x||e^{i\pi}e_1 = (11, 0, 0, 0)$. Using $v_1 = x - y = (-16, -4, 8, 4)^T$, we find that

$$Q_{v_1} = \frac{1}{11} \begin{bmatrix} -5 & -4 & 8 & 4 \\ -4 & 10 & 2 & 1 \\ 8 & 2 & 7 & -2 \\ 4 & 1 & -2 & 10 \end{bmatrix}, \quad Q_{v_1}A = \frac{1}{11} \begin{bmatrix} 121 & 41 & 21 & -69 \\ 0 & 13 & 63 & 112 \\ 0 & 29 & -82 & 40 \\ 0 & 97 & -19 & -24 \end{bmatrix}.$$

In the second step, we look at the lower-right 3×3 block and take $x = (13/11, 29/11, 97/11)^T$, $||x|| = 9.27941$, $x(1) = (13/11)e^{i0}$, so $y = (-9.27941, 0, 0)^T$. Taking

$v_2' = x - y$ gives

$$v_2 = (0, 10.4612, 2.63636, 8.81818)^T,$$

$$Q_{v_2} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & -0.127 & -0.284 & -0.950 \\ 0 & -0.284 & 0.928 & -0.239 \\ 0 & -0.950 & -0.239 & 0.199 \end{bmatrix}, Q_{v_2} Q_{v_1} A = \begin{bmatrix} 11 & 3.727 & 1.909 & -6.273 \\ 0 & -9.279 & 3.030 & -0.257 \\ 0 & 0 & -8.134 & 1.006 \\ 0 & 0 & -4.001 & -10.981 \end{bmatrix}.$$

In the third step, $x = (-8.134, -4.001)^T$, $\|x\| = 9.065$, $y = (9.065, 0)^T$, and $v_3 = (0, 0, -17.199, -4.001)^T$, so

$$Q_{v_3} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & -0.897 & -0.441 \\ 0 & 0 & -0.441 & 0.897 \end{bmatrix}, Q_{v_3} Q_{v_2} Q_{v_1} A = \begin{bmatrix} 11 & 3.727 & 1.909 & -6.273 \\ 0 & -9.279 & 3.030 & -0.257 \\ 0 & 0 & 9.065 & 3.944 \\ 0 & 0 & 0 & -10.297 \end{bmatrix}.$$

Finally, let D be diagonal with diagonal entries $1, -1, 1, -1$, and let $Q = Q_{v_1} Q_{v_2} Q_{v_3} D^{-1}$. Then $A = QR$ holds with

$$Q = \begin{bmatrix} -0.455 & 0.506 & -0.728 & 0.086 \\ -0.364 & 0.254 & 0.492 & 0.749 \\ 0.727 & 0.031 & -0.363 & 0.581 \\ 0.363 & 0.824 & 0.309 & -0.306 \end{bmatrix}, R = \begin{bmatrix} 11 & 3.727 & 1.909 & -6.273 \\ 0 & 9.279 & -3.030 & 0.257 \\ 0 & 0 & 9.065 & 3.944 \\ 0 & 0 & 0 & 10.297 \end{bmatrix}. \quad (9.2)$$

9.6 LU Factorization

In the next few sections, we discuss decompositions of a square matrix into the product of a lower-triangular matrix and an upper-triangular matrix. In order to assure the existence and uniqueness of such a decomposition, we must impose certain hypotheses on the given matrix A . Specifically, for any field F and $A \in M_n(F)$, let $A[k]$ denote the matrix in $M_k(F)$ consisting of the first k rows and columns of A . We will show that *for all $A \in M_n(F)$, there exist a lower-unitriangular matrix $L \in M_n(F)$ and an invertible upper-triangular matrix $U \in M_n(F)$ such that $A = LU$ iff for all $k \in [n]$, $\det(A[k]) \neq 0_F$. When L and U exist, they are unique.*

On one hand, assume $A = LU$ for some matrices L and U with the stated properties. In particular, $0 \neq \det(U) = \prod_{i=1}^n U(i, i)$. Fix $k \in [n]$. By writing $A = LU$ as a partitioned matrix product

$$\left[\begin{array}{c|c} A[k] & B \\ \hline C & D \end{array} \right] = \left[\begin{array}{c|c} L[k] & 0 \\ \hline M & N \end{array} \right] \left[\begin{array}{c|c} U[k] & V \\ \hline 0 & W \end{array} \right], \quad (9.3)$$

we see that $A[k] = L[k]U[k]$, where $L[k]$ is lower-unitriangular and $U[k]$ is upper-triangular. Taking determinants gives $\det(A[k]) = \det(L[k]) \det(U[k]) = \prod_{i=1}^k L(i, i) \prod_{i=1}^k U(i, i) \neq 0$. So the determinant condition on A is necessary for the LU factorization to exist.

On the other hand, to prove uniqueness, assume $A = LU = L_1 U_1$ where both L, U and L_1, U_1 have the properties stated in the theorem. Since L_1 and U are invertible, we can write $L_1^{-1}L = U_1U^{-1}$. The matrix $L_1^{-1}L$ is lower-unitriangular, and the matrix U_1U^{-1} is upper-triangular. These two matrices are equal, so $L_1^{-1}L$ is upper-triangular and lower-unitriangular. The only such matrix is I_n , which implies $L = L_1$ and $U = U_1$.

Finally, we assume $\det(A[k]) \neq 0$ for $1 \leq k \leq n$ and prove existence of the factorization

$A = LU$ by the following recursive computation. Define $L(i, i) = 1_F$ for $i \in [n]$ and $L(i, j) = U(j, i) = 0_F$ for $1 \leq i < j \leq n$. We determine the remaining entries of L one column at a time, and simultaneously determine the remaining entries of U one row at a time. Fix $k \in [n]$, and assume inductively that we have already determined all entries in the first $k - 1$ columns of L and all entries in the first $k - 1$ rows of U . Assume further that $A(i, j) = (LU)(i, j)$ holds for all $i, j \in [n]$ such that $i \leq k - 1$ or $j \leq k - 1$. (For the base case $k = 1$, these assumptions hold vacuously.)

To continue, we first recover $U(k, k)$ as follows. Comparing the k, k -entries of A and LU , we see that these entries will be equal iff

$$A(k, k) = \sum_{r=1}^n L(k, r)U(r, k) = \sum_{r=1}^{k-1} L(k, r)U(r, k) + L(k, k)U(k, k). \quad (9.4)$$

By hypothesis, every entry in the sum from $r = 1$ to $k - 1$ is already known. Since we chose $L(k, k) = 1$, the preceding equation holds iff we set

$$U(k, k) = A(k, k) - \sum_{r=1}^{k-1} L(k, r)U(r, k). \quad (9.5)$$

At this point, we know (among other things) that $A(i, j) = (LU)(i, j)$ for all $i, j \in [k]$ and that the zero blocks displayed in (9.3) are present in L and U . So the equation $A[k] = L[k]U[k]$ is valid. Taking determinants shows that $0 \neq \det(A[k]) = \det(L[k])\det(U[k]) = \prod_{i=1}^k U(i, i)$. Since the field F has no zero divisors, we conclude that $U(k, k) \neq 0_F$.

The next step is to compute the unknown entries in column k of L , which are the entries $L(i, k)$ as i ranges from $k + 1$ to n . Comparing the i, k -entries of A and LU for each i in this range, we see that these entries will agree iff

$$A(i, k) = \sum_{r=1}^n L(i, r)U(r, k) = \sum_{r=1}^{k-1} L(i, r)U(r, k) + L(i, k)U(k, k).$$

The entries $L(i, r)$ and $U(r, k)$ for $1 \leq r \leq k - 1$ are already known, and we have seen that $U(k, k) \neq 0_F$. Thus, $A(i, k) = (LU)(i, k)$ will hold iff we set

$$L(i, k) = U(k, k)^{-1} \left[A(i, k) - \sum_{r=1}^{k-1} L(i, r)U(r, k) \right]. \quad (9.6)$$

We now know all of column k of L , and we have ensured that $A(i, j) = (LU)(i, j)$ whenever $i < k$ or $j \leq k$.

The last step is to compute the unknown entries in row k of U , which are the entries $U(k, j)$ as j ranges from $k + 1$ to n . Comparing the k, j -entries of A and LU for each j in this range, we see that these entries will agree iff

$$A(k, j) = \sum_{r=1}^n L(k, r)U(r, j) = \sum_{r=1}^{k-1} L(k, r)U(r, j) + L(k, k)U(k, j).$$

The entries $L(k, r)$ and $U(r, j)$ for $1 \leq r \leq k - 1$ are already known, and we chose $L(k, k) = 1_F$. So $A(k, j) = (LU)(k, j)$ will hold iff we set

$$U(k, j) = A(k, j) - \sum_{r=1}^{k-1} L(k, r)U(r, j). \quad (9.7)$$

We now know all of row k of U , and we have ensured that $A(i, j) = (LU)(i, j)$ whenever $i \leq k$ or $j \leq k$. This completes step k of the induction, and the existence proof is now finished.

The factorization $A = LU$ just proved, where L is lower-unitriangular and U is upper-triangular and invertible, is sometimes called the *Doolittle factorization* of A . Some variants of this factorization can be found by changing the conditions on the diagonal entries of L and U . For example, assuming $\det(A[k]) \neq 0$ for $1 \leq k \leq n$, there exist unique L, D, U in $M_n(F)$ such that L is lower-unitriangular, D is diagonal and invertible, and U is upper-unitriangular. With the same assumption on A , there exist unique $L, U \in M_n(F)$ such that L is lower-triangular and invertible, and U is upper-unitriangular; this is called the *Croft factorization* of A . In the field $F = \mathbb{C}$ (or any field where every element has a square root), we can write $A = LU$ with $L(k, k) = U(k, k)$ for all k by setting $L(k, k) = U(k, k)$ in (9.4) and solving for $U(k, k)$. We ask the reader to derive these variant factorizations in Exercise 34.

9.7 Example of the LU Factorization

We illustrate the computation of an LU factorization using the matrix

$$A = \begin{bmatrix} 3 & 1 & -5 & 2 \\ -9 & -1 & 14 & -5 \\ 15 & 9 & -23 & 12 \\ 3 & 9 & -9 & 7 \end{bmatrix},$$

which does satisfy the determinant condition $\det(A[k]) \neq 0$ for $1 \leq k \leq 4$. Initially, the unknown factors L and U look like this:

$$L = \begin{bmatrix} 1 & 0 & 0 & 0 \\ a & 1 & 0 & 0 \\ b & d & 1 & 0 \\ c & e & f & 1 \end{bmatrix}, \quad U = \begin{bmatrix} q & r & s & t \\ 0 & u & v & w \\ 0 & 0 & x & y \\ 0 & 0 & 0 & z \end{bmatrix}.$$

Comparing the entries of the matrix equation $A = LU$ leads to the following sixteen scalar equations in the unknowns a, b, c, \dots, z :

$$\begin{array}{llll} 3 = q & 1 = r & -5 = s & 2 = t \\ -9 = aq & -1 = ar + u & 14 = as + v & -5 = at + w \\ 15 = bq & 9 = br + du & -23 = bs + dv + x & 12 = bt + dw + y \\ 3 = cq & 9 = cr + eu & -9 = cs + ev + fx & 7 = ct + ew + fy + z \end{array}$$

We see $q = 3$, then (working down the first column) $a = -3$, $b = 5$, $c = 1$, and (looking at the first row) $r = 1$, $s = -5$, $t = 2$. From the 2, 2-entry, we now see $u = -1 - ar = 2$. Working down the second column, $9 = 5 + 2d$ gives $d = 2$, and $9 = 1 + 2e$ gives $e = 4$. Moving along the second row, $14 = 15 + v$ gives $v = -1$, and $-5 = -6 + w$ gives $w = 1$. The 3, 3-entry gives $-23 = -25 - 2 + x$, so $x = 4$. Continuing similarly, we find $f = 0$, $y = 0$, and $z = 1$. So, the Doolittle version of the LU factorization is

$$A = \begin{bmatrix} 3 & 1 & -5 & 2 \\ -9 & -1 & 14 & -5 \\ 15 & 9 & -23 & 12 \\ 3 & 9 & -9 & 7 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ -3 & 1 & 0 & 0 \\ 5 & 2 & 1 & 0 \\ 1 & 4 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 3 & 1 & -5 & 2 \\ 0 & 2 & -1 & 1 \\ 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

We can force U to be unitriangular by factoring out an appropriate diagonal matrix:

$$A = \begin{bmatrix} 1 & 0 & 0 & 0 \\ -3 & 1 & 0 & 0 \\ 5 & 2 & 1 & 0 \\ 1 & 4 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 3 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1/3 & -5/3 & 2/3 \\ 0 & 1 & -1/2 & 1/2 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

Then we can absorb the diagonal matrix into L to obtain the Crout factorization of A :

$$A = \begin{bmatrix} 3 & 0 & 0 & 0 \\ -9 & 2 & 0 & 0 \\ 15 & 4 & 4 & 0 \\ 3 & 8 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1/3 & -5/3 & 2/3 \\ 0 & 1 & -1/2 & 1/2 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

9.8 LU Factorizations and Gaussian Elimination

There is a close relationship between LU factorizations and Gaussian elimination. To see how this arises, let us perform elementary row operations on the matrix A in §9.7 to eliminate the zeroes below the diagonal of A . The first step is to remove the -9 in the $2, 1$ -entry by adding 3 times row 1 of A to row 2 of A . Recall from §4.9 that this row operation can be achieved by multiplying A on the left by the *elementary matrix* E_1 that is obtained from I_4 by adding 3 times row 1 to row 2. The new matrix is

$$E_1 A = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 3 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 3 & 1 & -5 & 2 \\ -9 & -1 & 14 & -5 \\ 15 & 9 & -23 & 12 \\ 3 & 9 & -9 & 7 \end{bmatrix} = \begin{bmatrix} 3 & 1 & -5 & 2 \\ 0 & 2 & -1 & 1 \\ 15 & 9 & -23 & 12 \\ 3 & 9 & -9 & 7 \end{bmatrix}.$$

We continue by adding -5 times row 1 to row 3 using an elementary matrix E_2 , which produces a new matrix

$$E_2(E_1 A) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ -5 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 3 & 1 & -5 & 2 \\ 0 & 2 & -1 & 1 \\ 15 & 9 & -23 & 12 \\ 3 & 9 & -9 & 7 \end{bmatrix} = \begin{bmatrix} 3 & 1 & -5 & 2 \\ 0 & 2 & -1 & 1 \\ 0 & 4 & 2 & 2 \\ 3 & 9 & -9 & 7 \end{bmatrix}.$$

Next we add -1 times row 1 to row 4, obtaining

$$E_3(E_2 E_1 A) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ -1 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 3 & 1 & -5 & 2 \\ 0 & 2 & -1 & 1 \\ 0 & 4 & 2 & 2 \\ 3 & 9 & -9 & 7 \end{bmatrix} = \begin{bmatrix} 3 & 1 & -5 & 2 \\ 0 & 2 & -1 & 1 \\ 0 & 4 & 2 & 2 \\ 0 & 8 & -4 & 5 \end{bmatrix}.$$

Moving to column 2, we add -2 times row 2 to row 3:

$$E_4(E_3 E_2 E_1 A) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & -2 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 3 & 1 & -5 & 2 \\ 0 & 2 & -1 & 1 \\ 0 & 4 & 2 & 2 \\ 0 & 8 & -4 & 5 \end{bmatrix} = \begin{bmatrix} 3 & 1 & -5 & 2 \\ 0 & 2 & -1 & 1 \\ 0 & 0 & 4 & 0 \\ 0 & 8 & -4 & 5 \end{bmatrix}.$$

Then we add -4 times row 2 to row 4:

$$E_5(E_4 E_3 E_2 E_1 A) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & -4 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 3 & 1 & -5 & 2 \\ 0 & 2 & -1 & 1 \\ 0 & 0 & 4 & 0 \\ 0 & 8 & -4 & 5 \end{bmatrix} = \begin{bmatrix} 3 & 1 & -5 & 2 \\ 0 & 2 & -1 & 1 \\ 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

There is already a zero in the 4,3-position, so we are finished.

Observe that the final matrix $(E_5 E_4 E_3 E_2 E_1)A$ is precisely the upper-triangular matrix U found in §9.7. Solving for A , we find that $A = LU$ where $L = E_1^{-1} E_2^{-1} E_3^{-1} E_4^{-1} E_5^{-1}$. The inverse of each elementary matrix E_i is the lower-triangular matrix obtained by negating the unique nonzero entry of E_i not on the diagonal. Using this observation, one readily confirms that L as defined here agrees with the L found in §9.7. Note that for $i > j$, $L(i,j)$ is the negative of the multiple of row j that was added to row i to make the i,j -entry of A become zero.

Now consider what happens when we perform Gaussian elimination on a general matrix $A \in M_n(F)$ satisfying $\det(A[k]) \neq 0$ for $1 \leq k \leq n$. The algorithm proceeds as follows, assuming no divisions by zero occur (we justify this assumption later):

```

for j=1 to n-1 do
    for i=j+1 to n do
        m(i,j) = A(i,j)/A(j,j);
        replace row i of A by row i minus m(i,j) times row j;
    end for;
end for.

```

To describe this algorithm using elementary matrices, let $E(i,j,c)$ be the elementary matrix with i,j -entry equal to c , diagonal entries equal to 1, and other entries equal to zero (for $i \neq j$ in $[n]$ and $c \in F$). For each column index j in the outer loop, the inner loop modifies the current version of the matrix A by left-multiplying by a product of elementary matrices

$$C_j = E(n,j,-m(n,j))E(n-1,j,-m(n-1,j))\cdots E(j+1,j,-m(j+1,j)), \quad (9.8)$$

which create zeroes in column j below the diagonal. The net effect of the algorithm is to replace A by the matrix

$$C_{n-1} \cdots C_2 C_1 A = U,$$

which is an upper-triangular matrix. Solving for A , we can write $A = LU$ with

$$L = C_1^{-1} C_2^{-1} \cdots C_{n-1}^{-1} I_n.$$

Observe that

$$C_j^{-1} = E(j+1,j,m(j+1,j))\cdots E(n-1,j,m(n-1,j))E(n,j,m(n,j)).$$

To find the entries of L , we can start with I_n and work from right to left, applying the elementary row operations encoded by each elementary matrix $E(i,j,m(i,j))$ as j goes from $n-1$ down to 1 and (for each j) i goes from n down to $j+1$. Because the elementary operations are applied in this particular order, one sees that the final result is a lower-unitriangular matrix L with i,j -entry $L(i,j) = m(i,j)$ for all $i > j$ in $[n]$. To summarize, *in the Doolittle factorization $A = LU$, the entries of L encode the multipliers needed when we use Gaussian elimination to reduce A to the upper-triangular matrix U .*

To complete our analysis, we need only justify the assumption that the computation $m(i,j) = A(i,j)/A(j,j)$ in the algorithm will never involve a division by zero. We show this by induction on the outer loop variable j . Fix j with $1 \leq j < n$, and assume no divisions by zero occurred in the first $j-1$ iterations of the outer loop. At this point, the algorithm has transformed the original matrix A into a new matrix

$$C_{j-1} \cdots C_2 C_1 A = U_{j-1},$$

where $L_{j-1} = C_{j-1} \cdots C_2 C_1$ is lower-unitriangular, and U_{j-1} has zeroes below the diagonal in the first $j-1$ columns. Write $L_{j-1}A = U_{j-1}$ as a partitioned matrix product

$$\left[\begin{array}{c|c} L_{j-1}[j] & 0 \\ \hline * & * \end{array} \right] \left[\begin{array}{c|c} A[j] & * \\ \hline * & * \end{array} \right] = \left[\begin{array}{c|c} U_{j-1}[j] & * \\ \hline * & * \end{array} \right], \quad (9.9)$$

where each upper-left block is $j \times j$. We see that $L_{j-1}[j]A[j] = U_{j-1}[j]$. Taking determinants and noting that $U_{j-1}[j]$ is upper-triangular, we get

$$0 \neq \det(A[j]) = \det(L_{j-1}[j]) \det(A[j]) = \det(U_{j-1}[j]) = \prod_{k=1}^j U_{j-1}(k, k).$$

So the value of $A(j, j)$ used in the j 'th iteration of the outer loop, namely $U_{j-1}(j, j)$, is nonzero as claimed.

9.9 Permuted LU Factorizations

Our next goal is to generalize the analysis in the preceding section to the case where $A \in M_n(F)$ is an invertible matrix, but A does not satisfy the condition $\det(A[k]) \neq 0$ for all $k \in [n]$. Such an A cannot have an LU factorization, but we will prove there is a permutation matrix $P \in M_n(F)$ such that PA has an LU factorization. Recall that a *permutation matrix* P has exactly one 1 in each row and column, with all other entries of P equal to zero. For $F = \mathbb{C}$, the columns of P are orthonormal (relative to the standard inner product on \mathbb{C}^n), so that $P^*P = P^T P = I_n$. It follows that P is invertible with $P^{-1} = P^T$; the same result holds for any field F (Exercise 36). By §4.8, PA is obtained from A by reordering the rows of A . Specifically, for each i, j with $P(i, j) = 1$, row i of PA is row j of A .

We now prove that *for any invertible $A \in M_n(F)$, there exist a permutation matrix $P \in M_n(F)$, a lower-unitriangular matrix $L \in M_n(F)$, and an upper-triangular invertible matrix $U \in M_n(F)$ with $PA = LU$, or equivalently $A = P^T LU$* . The idea of the proof is to reduce A to the upper-triangular matrix U via Gaussian elimination. The algorithm used before still works, with one modification. When we try to compute the multipliers $m(i, j)$ to process column j of the current matrix, we may find that $A(j, j)$ is zero. Recall that “the current matrix” in the j 'th iteration of the outer loop can be written $U_{j-1} = L_{j-1}A$, where A is the original input matrix, L_{j-1} is lower-unitriangular, and U_{j-1} has zeroes below the diagonal in the first $j-1$ columns. Since A and L_{j-1} are both invertible, U_{j-1} is invertible. Writing U_{j-1} in partitioned form as

$$U_{j-1} = \left[\begin{array}{c|c} U_{j-1}[j-1] & * \\ \hline 0 & * \end{array} \right], \quad (9.10)$$

it follows by taking determinants that for some $k \geq j$, $U_{j-1}(k, j) \neq 0$. Choose k minimal with this property and interchange rows j and k at the start of the j 'th iteration of the outer loop. The rest of the algorithm proceeds as before.

In terms of elementary matrices, the modified elimination algorithm yields a factorization

$$C_{n-1}P_{n-1} \cdots C_3P_3C_2P_2C_1P_1A = U, \quad (9.11)$$

where C_j has the form (9.8), and each P_j is an elementary permutation matrix obtained

from I_n by switching row j with some row $k_j \geq j$. To obtain the required factorization $PA = LU$, we need to move all the permutation matrices P_j to the right through all the matrices C_r . This manipulation requires some care, since P_j and C_r do not commute in general.

Consider a product $P_j C_r$ where $1 \leq r < j \leq n$. Each C_r is a product of elementary matrices of the form $E(s, r, m)$ with $s > r$ and $m \in F$. Suppose P_j is obtained by switching rows j and $k \geq j$ of I_n . We claim that $P_j E(s, r, m) = E(s', r, m)P_j$, where $s' = k$ if $s = j$, $s' = j$ if $s = k$, and $s' = s$ otherwise. We verify that $P_j E(j, r, m) = E(k, r, m)P_j$, leaving the other assertions in the claim as exercises. It suffices to show that $P_j E(j, r, m)B = E(k, r, m)P_j B$ for any matrix $B \in M_n(F)$. The left side is the matrix obtained from B by first adding m times row r to row j , then switching rows j and k . The right side is the matrix obtained from B by first switching rows j and k , then adding m times row r to row k . Since r is different from j and k , the effect of these two operations on B is the same. The rest of the claim is verified similarly.

By using the claim repeatedly to move P_j to the right past each factor in C_r , we see that for all $r < j$, $P_j C_r = C'_r P_j$ where C'_r is also a product of elementary matrices $E(s', r, m)$ with s' ranging from $r+1$ to n in some order (cf. (9.8)). Using this in (9.11), we can first write $P_2 C_1 = C'_1 P_2$ and continue by rewriting $C_3 P_3 C_2 C'_1 P_2 P_1$ as $C_3 C'_2 C''_1 P_3 P_2 P_1$. Ultimately, we will reach a new factorization

$$\tilde{C}_{n-1} \cdots \tilde{C}_2 \tilde{C}_1 (P_{n-1} \cdots P_2 P_1) A = U,$$

where each \tilde{C}_j has the form (9.8) (with factors appearing in a different order). The product of all P_j 's is a permutation matrix P . We now have $PA = LU$, where $L = \tilde{C}_1^{-1} \cdots \tilde{C}_{n-1}^{-1}$ is lower-unitriangular.

Finally, we extend our results to the case of a general square matrix A (not necessarily invertible). We prove that *for any $A \in M_n(F)$, there exist permutation matrices $P, Q \in M_n(F)$, a lower-unitriangular $L \in M_n(F)$, and an upper-triangular $U \in M_n(F)$ with $\text{rank}(A) = \text{rank}(U)$, such that $PAQ = LU$.* Examining the preceding proof, we see that the only place we needed invertibility of A was when we used (9.10) to show $U_{j-1}(k, j) \neq 0$ for some $k \geq j$. Consider what happens when our algorithm reduces a general matrix A to the matrix U_{j-1} given in (9.10).

Case 1: there is a nonzero entry somewhere in the lower-right block shown in (9.10), say $U_{j-1}(k, p) \neq 0$ for some $k, p \geq j$. We can multiply the current matrix *on the right* by the permutation matrix Q_j that switches *columns* j and p , and then multiply on the left (as before) by the permutation matrix P_j that switches rows j and k , to create a matrix with a nonzero entry in the j, j -position. The algorithm proceeds to create zeroes below this entry, as before.

Case 2: all entries in the lower-right block of (9.10) are zero. Then the algorithm terminates at this point with a factorization

$$C_{j-1} P_{j-1} \cdots C_1 P_1 A Q_1 Q_2 \cdots Q_{j-1} = U_{j-1} = U,$$

where U is upper-triangular since $U_{j-1}[j-1]$ is upper-triangular and rows j through n of U_{j-1} contain all zeroes. Moving the matrices P_1, \dots, P_{j-1} to the right (as in the invertible case) and rearranging, we obtain the required factorization $PAQ = LU$. Since P, Q , and L are all invertible, it readily follows that $\text{rank}(A) = \text{rank}(U)$.

9.10 Cholesky Factorization

This section proves the following *Cholesky factorization theorem* for positive semidefinite matrices: *for every positive semidefinite matrix $A \in M_n(\mathbb{C})$, there exists a lower-triangular matrix L with nonnegative diagonal entries such that $A = LL^*$. If A is positive definite (i.e., invertible), then L is unique.*

To begin the proof, recall from §7.11 that a positive semidefinite matrix $A \in M_n(\mathbb{C})$ has a unique positive semidefinite square root, which is a positive semidefinite matrix $B \in M_n(\mathbb{C})$ such that $A = B^2$. Since B must be Hermitian, we have $A = B^*B$. Now, B has a QR factorization $B = QR$, where $Q \in M_n(\mathbb{C})$ is unitary and $R \in M_n(\mathbb{C})$ is upper-triangular with nonnegative diagonal entries. Let $L = R^*$, which is lower-triangular with nonnegative diagonal entries. Since Q is unitary, we have $A = B^*B = (QR)^*QR = R^*Q^*QR = R^*R = LL^*$, as needed.

To prove uniqueness when A is invertible, suppose $A = LL^* = MM^*$ where both $L, M \in M_n(\mathbb{C})$ are lower-triangular with nonnegative diagonal entries. Taking determinants, we see that L and M are both invertible. Then $M^{-1}L = M^*(L^*)^{-1}$, where the left side is a lower-triangular matrix and the right side is an upper-triangular matrix. Since the two sides are equal, $M^{-1}L$ must be a diagonal matrix D with strictly positive diagonal entries. Then $L = MD$, so $MM^* = LL^* = MDD^*M^* = MD^2M^*$. Left-multiplying by M^{-1} and right-multiplying by $(M^*)^{-1}$ gives $I = D^2$. Since all diagonal entries of D are positive real numbers, this forces $D = I$ and $M = L$, completing the uniqueness proof.

One approach to computing the Cholesky factorization of a positive definite matrix A is to observe that $A = LL^*$ is one of the variations of the LU factorization from §9.6, in which L is lower-triangular, $L^* = U$ is upper-triangular, and $L(k, k) = L^*(k, k)$ for all $k \in [n]$. We can recover the columns of L (and hence the rows of $U = L^*$) one at a time using the formulas in §9.6. Specifically, for $k = 1$ to n , we first calculate

$$L(k, k) = \sqrt{A(k, k) - \sum_{r=1}^{k-1} L(k, r)L^*(r, k)}$$

and then compute

$$L(i, k) = L(k, k)^{-1} \left[A(i, k) - \sum_{r=1}^{k-1} L(i, r)L^*(r, k) \right]$$

for $i = k+1, \dots, n$. Since we proved above that there is a factorization of the form $A = LL^*$, the uniqueness of the LU factorization ensures that the U computed by the LU algorithm must equal L^* . In particular, there is no need to use (9.7) to solve for the entries of U in this setting.

For example, consider $A = \begin{bmatrix} 4 & 2 & -1 \\ 2 & 5 & 1 \\ -1 & 1 & 3 \end{bmatrix}$. By the determinant test in §7.15, A is positive definite. We seek a factorization $A = LL^*$, say

$$\begin{bmatrix} 4 & 2 & -1 \\ 2 & 5 & 1 \\ -1 & 1 & 3 \end{bmatrix} = \begin{bmatrix} r & 0 & 0 \\ s & u & 0 \\ t & v & w \end{bmatrix} \begin{bmatrix} r & s & t \\ 0 & u & v \\ 0 & 0 & w \end{bmatrix}.$$

(Since A is real and positive definite, L will be real and $L^* = L^T$.) Comparing 1, 1-entries,

we find $r^2 = 4$ and $r = 2$. Working down the first column, $sr = 2$ and $tr = -1$, so $s = 1$ and $t = -1/2$. Moving to the 2,2-entry, $s^2 + u^2 = 5$ becomes $1 + u^2 = 5$, so $u = 2$. Looking at the 3,2-entry, $1 = st + vu$ implies $1 = -1/2 + 2v$, so $v = 3/4$. Finally, $t^2 + v^2 + w^2 = 3$ gives $1/4 + 9/16 + w^2 = 3$, so $w = \sqrt{35/16} = \sqrt{35}/4$. We can check the computation by confirming that

$$\begin{bmatrix} 4 & 2 & -1 \\ 2 & 5 & 1 \\ -1 & 1 & 3 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 \\ 1 & 2 & 0 \\ -1/2 & 3/4 & \sqrt{35}/4 \end{bmatrix} \begin{bmatrix} 2 & 1 & -1/2 \\ 0 & 2 & 3/4 \\ 0 & 0 & \sqrt{35}/4 \end{bmatrix}.$$

9.11 Least Squares Approximation

This section describes an application of the Cholesky factorization to the solution of a *least squares approximation problem*. We want to solve a linear system of the form $Ax = b$, where $A \in M_{m,n}(\mathbb{C})$ is a given matrix with $m \geq n$, $b \in \mathbb{C}^m$ is a given column vector, and $x \in \mathbb{C}^n$ is unknown. This linear system typically has many more equations than unknowns, so we would not expect there to be any exact solutions. Instead, we seek an *approximate* solution x that minimizes the length of the *error vector* $Ax - b$. The squared length of this error vector is $\sum_{i=1}^m |(Ax)_i - b_i|^2$, which is the sum of the squares of the errors in each individual equation.

We assume that $\text{rank}(A) = n$, so that A has n linearly independent columns. Under this assumption, we now show that *there exists a unique $x \in \mathbb{C}^n$ minimizing $\|Ax - b\|^2 = \langle Ax - b, Ax - b \rangle$, and this x is the unique solution to the square linear system $(A^*A)x = A^*b$.* The linear equations in this new system are called *normal equations* for x , and x is called the *least squares approximate solution* to $Ax = b$.

We first check that $A^*A \in M_n(\mathbb{C})$ is invertible. Fix $z \in \mathbb{C}^n$ with $A^*Az = 0$; it suffices to show $z = 0$ (see Table 4.1). Note that $\|Az\|^2 = \langle Az, Az \rangle = (Az)^*(Az) = z^*A^*Az = 0$, so $Az = 0$. If z were not zero, then writing $Az = 0$ as a linear combination of the n columns of A (as explained in §4.7) would show that these columns were linearly dependent, contrary to hypothesis. So $z = 0$, A^*A is invertible, and the normal equations have a unique solution $x = (A^*A)^{-1}(A^*b)$.

Now let $y \in \mathbb{C}^n$ be arbitrary. We claim the vectors $Ax - b$ and $A(y - x)$ are orthogonal. For, $\langle Ax - b, A(y - x) \rangle = (y - x)^*A^*(Ax - b) = (y - x)^*(A^*Ax - A^*b) = 0$ since x solves the normal equations. Using the Pythagorean theorem, we compute

$$\|Ay - b\|^2 = \|A(y - x) + (Ax - b)\|^2 = \|A(y - x)\|^2 + \|Ax - b\|^2 \geq \|Ax - b\|^2.$$

This inequality shows that among all vectors in \mathbb{C}^n , the x chosen above minimizes $\|Ax - b\|^2$. Furthermore, equality can only hold if $\|A(y - x)\|^2 = 0$, which implies $A(y - x) = 0$, hence (as seen in the last paragraph) $y - x = 0$ and $y = x$. So the vector x minimizing the sum of the squares of the errors is unique, as needed.

To see how the Cholesky factorization can be used here, note that A^*A is positive definite, since for all nonzero $v \in \mathbb{C}^n$, $Av \neq 0$ and so $v^*A^*Av = \langle Av, Av \rangle = \|Av\|^2 \in \mathbb{R}^+$. We can therefore write $A^*A = LL^*$ for a unique lower-triangular $L \in M_n(\mathbb{C})$. The normal equations characterizing the least squares solution x become $LL^*x = A^*b$. If we have found L , we can quickly solve the normal equations by first solving $Lw = A^*b$ using forward substitution, then solving $L^*x = w$ using backward substitution. We must point out, however, that solving least squares problems via the normal equations can lead to issues of numerical stability. Exercise 48 describes another approach to finding x based on a QR-factorization of A .

9.12 Singular Value Decomposition

We conclude the chapter by proving the *singular value decomposition theorem* for complex matrices: *given any $A \in M_{m,n}(\mathbb{C})$, there exist a unitary matrix $U \in M_n(\mathbb{C})$, a unitary matrix $V \in M_m(\mathbb{C})$, and a diagonal matrix $D \in M_{m,n}(\mathbb{R})$ with diagonal entries $d_1 \geq d_2 \geq \dots \geq d_{\min(m,n)} \geq 0$, such that $A = VDU^*$. The numbers d_j are uniquely determined by A ; the nonzero d_j 's are the positive square roots of the nonzero eigenvalues of A^*A .* We call $d_1, \dots, d_{\min(m,n)}$ the *singular values* of A .

We proved the singular value decomposition for square matrices as a consequence of the polar decomposition in §7.13. This section gives an independent proof valid for rectangular matrices. Fix $A \in M_{m,n}(\mathbb{C})$. The matrix $A^*A \in M_n(\mathbb{C})$ is positive semidefinite, since for any $x \in \mathbb{C}^n$, $x^*A^*Ax = (Ax)^*(Ax) = \langle Ax, Ax \rangle \geq 0$. We saw in Chapter 7 that positive semidefinite matrices are normal, hence unitarily diagonalizable via an orthonormal basis of eigenvectors, and all their eigenvalues are nonnegative real numbers. Let the eigenvalues of A^*A (with multiplicity) be $c_1 \geq c_2 \geq \dots \geq c_n \geq 0$. Let $u_1, u_2, \dots, u_n \in \mathbb{C}^n$ be associated unit eigenvectors, so (u_1, \dots, u_n) is an orthonormal basis of \mathbb{C}^n and $A^*Au_j = c_j u_j$ for $1 \leq j \leq n$. Let $d_j = \sqrt{c_j} \geq 0$ for $1 \leq j \leq n$, and let r be the maximal index with $c_j > 0$. For $1 \leq j \leq r$, define $v_j = d_j^{-1}Au_j \in \mathbb{C}^m$. We claim (v_1, \dots, v_r) is an orthonormal list in \mathbb{C}^m . To see why, fix $i, j \in [r]$ and compute

$$\langle v_j, v_i \rangle = (d_i^{-1}Au_i)^*(d_j^{-1}Au_j) = d_i^{-1}d_j^{-1}u_i^*A^*Au_j = d_i^{-1}d_j^{-1}u_i^*(c_j u_j) = (d_j/d_i)\langle u_j, u_i \rangle.$$

Since the u_k 's are orthonormal, we get $\langle v_j, v_i \rangle = 0$ if $j \neq i$, and $\langle v_j, v_i \rangle = d_j/d_i = 1$ if $j = i$.

Using Gram–Schmidt orthonormalization, we can extend the orthonormal list (v_1, \dots, v_r) to an orthonormal basis $(v_1, \dots, v_r, \dots, v_m)$ of \mathbb{C}^m . Let $U \in M_n(\mathbb{C})$ be the matrix with columns u_1, \dots, u_n , and let $V \in M_m(\mathbb{C})$ be the matrix with columns v_1, \dots, v_m . Both U and V are unitary matrices, since they are square matrices with orthonormal columns (see §7.5). Let $D \in M_{m,n}(\mathbb{C})$ have diagonal entries $D(j,j) = d_j$ for $1 \leq j \leq \min(m,n)$, and all other entries zero. To prove that $A = VDU^*$, we prove the equivalent identity $AU = VD$ (recall $U^* = U^{-1}$ since U is unitary). Fix $j \in [n]$. On the left side, the j 'th column of AU is $A(U^{[j]}) = Au_j$. On the right side, the j 'th column of VD is $V(D^{[j]}) = D(j,j)V^{[j]} = d_j v_j$. If $j \leq r$, then $d_j v_j = Au_j$ holds by definition of v_j . If $r < j \leq n$, then $c_j = d_j = 0$. So $\|Au_j\|^2 = u_j^*A^*Au_j = u_j^*(c_j u_j) = 0$, hence $Au_j = 0 = d_j v_j$. We conclude $AU = VD$ since all columns agree.

To prove the uniqueness of the d_j 's, let $A = VDU^*$ be any factorization with the properties given in the theorem statement (not necessarily using the matrices U, D, V constructed above). Since V is unitary, $A^*A = (VDU^*)^*(VDU^*) = UD^*V^*VDU^* = UD^*DU^*$. Since U is unitary, A^*A and D^*D are similar and so have the same eigenvalues (counted with multiplicity). Letting the nonzero diagonal entries of D be d_1, \dots, d_r , one readily checks that D^*D has eigenvalues d_1^2, \dots, d_r^2 , together with $n - r$ eigenvalues equal to zero. So the nonzero diagonal entries of D are the positive square roots of the nonzero eigenvalues of A^*A , as claimed.

We remark that the matrices V and U appearing in the singular value decomposition are not unique in general. We can obtain different matrices U by picking different orthonormal bases for each eigenspace of A^*A when forming the list (u_1, \dots, u_n) . Each choice of U leads to a unique list (v_1, \dots, v_r) , but then we can get different matrices V by extending this list to an orthonormal basis of \mathbb{C}^m in different ways.

In terms of linear transformations, the singular value decomposition can be phrased as follows. *For any linear map $T : \mathbb{C}^n \rightarrow \mathbb{C}^m$, there exist an orthonormal basis (u_1, \dots, u_n) of \mathbb{C}^n , an orthonormal basis (v_1, \dots, v_m) of \mathbb{C}^m , and nonnegative numbers $d_1 \geq \dots \geq d_{\min(m,n)}$*

such that $T(u_j) = d_j v_j$ for $1 \leq j \leq \min(m, n)$, and $T(u_j) = 0$ for $m < j \leq n$. We can also view the factorization $A = VDU^*$ geometrically, by saying that the linear transformation $x \mapsto Ax$ acts as the composition of an isometry on \mathbb{C}^n encoded by the unitary matrix U^* , followed by a nonnegative rescaling of the coordinate axes encoded by D (which deletes the last $n - m$ coordinates when $n > m$, and adds $m - n$ zero coordinates when $m > n$), followed by an isometry on \mathbb{C}^m encoded by the unitary matrix V .

9.13 Summary

1. *Orthogonality and Orthonormality.* A list (u_1, \dots, u_k) in an inner product space V is *orthogonal* iff $\langle u_i, u_j \rangle = 0$ for all $i \neq j$ in $[k]$. This list is *orthonormal* iff it is orthogonal and $\langle u_i, u_i \rangle = 1$ for all $i \in [k]$. Orthogonal lists of nonzero vectors are linearly independent. The *Pythagorean theorem* asserts that $\|x + y\|^2 = \|x\|^2 + \|y\|^2$ for all orthogonal $x, y \in V$. A square matrix Q is unitary iff its columns are orthonormal iff $Q^*Q = I$ iff $Q^* = Q^{-1}$.
2. *Orthogonal Projections.* Let (u_1, \dots, u_k) be an orthonormal list in an inner product space V . Let U be the subspace spanned by this list. For any $v \in V$, there exists a unique $x \in U$ minimizing $\|v - x\|$, namely $x = \sum_{r=1}^k \langle v, u_r \rangle u_r$, and $v - x$ is orthogonal to every vector in U . We call x the *orthogonal projection of v onto U* . The vector v lies in U iff $v = x$ iff $v = \sum_{r=1}^k \langle v, u_r \rangle u_r$.
3. *Gram–Schmidt Orthonormalization.* Given a finite or countably infinite linearly independent sequence (v_1, v_2, \dots) in an inner product space V , the *Gram–Schmidt orthonormalization algorithm* explicitly computes an orthonormal sequence (u_1, u_2, \dots) such that (v_1, \dots, v_r) and (u_1, \dots, u_r) span the same subspace V_r for all $r \geq 1$. Having found u_1, \dots, u_{s-1} , we calculate $x_s = \sum_{r=1}^{s-1} \langle v_s, u_r \rangle u_r$ and $u_s = (v_s - x_s)/\|v_s - x_s\|$. When $\dim(V) < \infty$, we see that every subspace of V has an orthonormal basis, and every orthonormal list in V can be extended to an orthonormal basis of V . If linear dependencies exist among the v_j 's, the algorithm detects this by computing $x_s = v_s$.
4. *QR Factorization via Gram–Schmidt.* Given $A \in M_{n,k}(\mathbb{C})$ with k linearly independent columns, there exist unique matrices $Q \in M_{n,k}(\mathbb{C})$ and $R \in M_k(\mathbb{C})$ such that $A = QR$, Q has orthonormal columns, and R is upper-triangular with strictly positive diagonal entries. When A is real, Q and R are real. If A has column rank $r < k$, we can write $A = QR$ where $Q \in M_{n,r}(\mathbb{C})$ has orthonormal columns and $R \in M_{r,k}(\mathbb{C})$ is upper-triangular with nonnegative diagonal entries. The columns of Q are the orthonormal vectors obtained by applying Gram–Schmidt orthonormalization to the columns of A . The entries of R are given by $R(i, j) = \langle v_j, u_i \rangle$ for $i \leq j$, and $R(i, j) = 0$ for $i > j$.
5. *Householder Reflections.* Given a nonzero $v \in \mathbb{C}^n$, the *Householder matrix* is $Q_v = I_n - (2/\|v\|^2)vv^* \in M_n(\mathbb{C})$. These matrices are unitary and self-inverse ($Q_v^* = Q_v^{-1} = Q_v$), $Q_v v = -v$, and $Q_v w = w$ for all $w \in \mathbb{C}^n$ orthogonal to v . For any $x \neq y$ in \mathbb{C}^n , there exists $v \in \mathbb{C}^n$ with $Q_v x = y$ iff $\|x\| = \|y\|$ and $\langle x, y \rangle \in \mathbb{R}$. Specifically, when x and y satisfy these conditions, $Q_v x = y$ iff v is a nonzero scalar multiple of $x - y$.
6. *Householder's QR Algorithm.* The Householder algorithm reduces $A \in M_{n,k}(\mathbb{C})$ to an upper-triangular matrix $R \in M_{n,k}(\mathbb{C})$ by applying $s = \min(n - 1, k)$

Householder reflections to create zeroes one column at a time. At the j 'th step, one finds a Householder reflection sending the partial column vector $(A(j, j), A(j + 1, j), \dots, A(n, j))^T$ to a vector whose last $n - j$ entries are zero. The output of the algorithm is a factorization $A = QR$ where Q is the product of s Householder matrices and a diagonal unitary matrix.

7. *LU Factorizations.* For any field F and $A \in M_n(F)$, there exist a lower-unitriangular matrix $L \in M_n(F)$ and an invertible upper-triangular matrix $U \in M_n(F)$ such that $A = LU$ iff for all $k \in [n]$, $\det(A[k]) \neq 0$. When L and U exist, they are unique. This is *Doolittle's LU factorization*; we obtain *Crout's LU factorization* by making U unitriangular instead of L . We can also write $A = LDU$ with both L and U unitriangular and D diagonal.
8. *Recursive Formula for LU Factorizations.* When the Doolittle factorization $A = LU$ exists, we can find L and U recursively as follows:

```

set L(i,i)=1 and L(i,j)=0=U(j,i) for all i<j in [n];
for k=1 to n do
    U(k,k)=A(k,k)-sum(r=1 to k-1) L(k,r)U(r,k);
    for i=k+1 to n do
        L(i,k)=(1/U(k,k))[A(i,k)-sum(r=1 to k-1) L(i,r)U(r,k)];
    end for;
    for j=k+1 to n do
        U(k,j)=A(k,j)-sum(r=1 to k-1) L(k,r)U(r,j);
    end for;
end for.

```

9. *LU Factorization via Gaussian Elimination.* For $A \in M_n(F)$ with $\det(A[k]) \neq 0$ for $k \in [n]$, the Gaussian elimination algorithm proceeds as follows:

```

for j=1 to n-1 do
    for i=j+1 to n do
        m(i,j) = A(i,j)/A(j,j);
        replace row i of A by row i minus m(i,j) times row j;
    end for;
end for.

```

In the Doolittle factorization $A = LU$, U is the output of the elimination algorithm, and the entries $L(i, j)$ below the diagonal are the multipliers $m(i, j)$ used in the algorithm.

10. *Permuted LU Factorizations.* For any $A \in M_n(F)$, there exist permutation matrices $P, Q \in M_n(F)$, a lower-unitriangular matrix $L \in M_n(F)$, and an upper-triangular matrix $U \in M_n(F)$ with $PAQ = LU$ and $\text{rank}(A) = \text{rank}(U)$. When A is invertible, we may take $Q = I_n$. These factorizations can be achieved by using Gaussian elimination with row and column interchanges.
11. *Cholesky Factorization.* Every positive semidefinite matrix $A \in M_n(\mathbb{C})$ can be factored as $A = LL^*$ for some lower-triangular $L \in M_n(\mathbb{C})$ with nonnegative diagonal entries. When A is positive definite, L is unique.
12. *Least Squares Approximation.* Given $b \in \mathbb{C}^m$ and $A \in M_{m,n}(\mathbb{C})$ of rank n with $m \geq n$, there exists a unique $x \in \mathbb{C}^n$ minimizing $\|Ax - b\|^2$, namely the unique solution to the normal equations $A^*Ax = A^*b$.

13. *Singular Value Decomposition.* Any $A \in M_{m,n}(\mathbb{C})$ can be factored as $A = VDU^*$, where $V \in M_m(\mathbb{C})$ and $U \in M_n(\mathbb{C})$ are unitary, and $D \in M_{m,n}(\mathbb{R})$ is diagonal with diagonal entries $d_1 \geq d_2 \geq \dots \geq 0$. The d_j 's are called the singular values of A and are uniquely determined by A ; the nonzero d_j 's are the positive square roots of the nonzero eigenvalues of A^*A .
-

9.14 Exercises

1. Let $L = \begin{bmatrix} 5 & 0 & 0 & 0 \\ -3 & 2 & 0 & 0 \\ 7 & -1 & -4 & 0 \\ 8 & -3 & -5 & 2 \end{bmatrix}$, $U = \begin{bmatrix} 1 & 3 & 2 & -4 \\ 0 & 1 & 5 & -3 \\ 0 & 0 & 1 & 6 \\ 0 & 0 & 0 & 1 \end{bmatrix}$, and $b = \begin{bmatrix} 10 \\ 6 \\ -12 \\ 4 \end{bmatrix}$.

(a) Solve $Lx = b$ by forward substitution. (b) Solve $Ux = b$ by backward substitution. (c) Solve $Ax = b$, where $A = LU$.

2. Let $Q = \begin{bmatrix} 0.28 & 0 & 0 & 0.96 \\ 0 & 0.6 & 0.8 & 0 \\ -0.96 & 0 & 0 & 0.28 \\ 0 & 0.8 & -0.6 & 0 \end{bmatrix}$, $R = \begin{bmatrix} 2 & 3 & -1 & -1 \\ 0 & 5 & 3 & -1 \\ 0 & 0 & 2 & -2 \\ 0 & 0 & 0 & 2 \end{bmatrix}$, and

$$b = \begin{bmatrix} 1 \\ 3 \\ 2 \\ -2 \end{bmatrix}. \text{ (a) Confirm that } Q \text{ is unitary. (b) Solve } Ax = b \text{ efficiently, where } A = QR.$$

3. (a) Suppose $A = LU$, where $A, L, U \in M_n(\mathbb{C})$ with L lower-unitriangular and U upper-triangular. If we know L and U , how many multiplications and divisions in \mathbb{C} are required to solve $Ax = b$? (b) Suppose $A = QR$, where $A, Q, R \in M_n(\mathbb{C})$ with Q unitary and R upper-triangular. If we know Q and R , how many multiplications and divisions in \mathbb{C} are required to solve $Ax = b$?
4. Apply the Gram–Schmidt algorithm to $v_1 = (1, 1, 1)$, $v_2 = (1, 2, 4)$, $v_3 = (1, 3, 9)$ to compute an orthonormal basis of \mathbb{C}^3 .
5. Let V be the subspace of \mathbb{C}^4 spanned by $v_1 = (3, 1, 5, -1)$, $v_2 = (2, 0, 1, 3)$, and $v_3 = (1, -1, -1, -1)$. (a) Apply the Gram–Schmidt algorithm to the list (v_1, v_2, v_3) to compute an orthonormal basis of V . (b) Find the orthogonal projection of each vector onto V : $e_1; (-2, 2, 3, -2); (1, 2, 3, 4); (24, -2, 15, 7)$.
6. Let V be a complex inner product space. (a) Prove: for all $x, y \in V$,

$$\|x + y\|^2 + \|x - y\|^2 = 2\|x\|^2 + 2\|y\|^2.$$

(b) Give an example where $x, y \in V$ and $\|x + y\|^2 = \|x\|^2 + \|y\|^2$, but x and y are not orthogonal. (c) Prove: for all $x, y \in V$, if $\|x + y\|^2 = \|x\|^2 + \|y\|^2 = \|x + iy\|^2$ then x and y are orthogonal.

7. Let V be the real inner product space of continuous functions $f : [0, 1] \rightarrow \mathbb{R}$, with inner product $\langle f, g \rangle = \int_0^1 f(x)g(x) dx$ for $f, g \in V$. (a) Apply Gram–Schmidt orthonormalization to the list $(1, x, x^2, x^3)$ to obtain an orthonormal basis for the subspace W spanned by this list. (b) Find the orthogonal projection of $f(x) = e^x$ onto W and the minimum distance from f to W . (c) Repeat part (b) for $f(x) = \sin(2\pi x)$.

8. Repeat Exercise 7, but let V consist of continuous functions $f : [-1, 1] \rightarrow \mathbb{R}$, with $\langle f, g \rangle = \int_{-1}^1 f(x)g(x) dx$ for $f, g \in V$.
9. (a) Suppose u_1, \dots, u_k are orthonormal vectors and $c_1, \dots, c_k \in \mathbb{C}$. Prove $\|c_1u_1 + \dots + c_ku_k\|^2 = |c_1|^2 + \dots + |c_k|^2$. (b) How does the formula in (a) change if we only assume u_1, \dots, u_k are orthogonal vectors?
10. Let (u_1, \dots, u_k) be an orthonormal list in an inner product space V , and let $v \in V$. Prove that $\sum_{j=1}^k |\langle v, u_j \rangle|^2 \leq \|v\|^2$, with equality iff v is in the span of u_1, \dots, u_k .
11. Let $V = \mathbb{R}[x]$ be the real inner product space of polynomials, with inner product

$$\left\langle \sum_{i=0}^n a_i x^i, \sum_{j=0}^m b_j x^j \right\rangle = \sum_{i=0}^{\min(m,n)} a_i b_i \text{ for } a_i, b_j \in \mathbb{R}.$$

- (a) Suppose $(f_0, f_1, f_2, \dots, f_n, \dots)$ is any infinite list of elements of V with $\deg(f_n) = n$ for all $n \geq 0$. Show that applying Gram–Schmidt orthonormalization to this list produces the sequence $(\pm 1, \pm x, \pm x^2, \dots, \pm x^n, \dots)$. (b) Given a finite subset $S \subseteq \mathbb{N}$, what is the orthogonal projection of $g \in \mathbb{R}[x]$ onto the subspace spanned by $(x^j : j \in S)$?
12. Use the Gram–Schmidt algorithm to prove that any orthonormal list (u_1, \dots, u_k) in a finite-dimensional inner product space V can be extended to an orthonormal ordered basis of V .
13. (a) Show that if $R, R_1 \in M_k(\mathbb{C})$ are upper-triangular with positive real diagonal entries, then RR_1^{-1} is upper-triangular with positive real diagonal entries. (b) Show that the only unitary upper-triangular matrix in $M_k(\mathbb{C})$ with positive real diagonal entries is I_k . (c) Use (a) and (b) to prove the uniqueness of the factorization $A = QR$ in the case where $A \in M_k(\mathbb{C})$ is a square matrix.
14. For each matrix A , use Gram–Schmidt orthonormalization to find the QR factorizations described in §9.3. (a) $A = \begin{bmatrix} 1 & 2 \\ 1 & 1 \end{bmatrix}$; (b) $A = \begin{bmatrix} 3 & 4 \\ 4 & 3 \end{bmatrix}$; (c) $A = \begin{bmatrix} 2 & 1 & 3 \\ 3 & 1 & 5 \\ 6 & 2 & -1 \end{bmatrix}$; (d) $A = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 3 & 3 & 3 & 3 \\ 2 & 0 & 2 & 0 \end{bmatrix}$; (e) $A = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 3 & 2 & 4 \\ 1 & 4 & 2 & 3 \\ 2 & 3 & 1 & 4 \end{bmatrix}$.
15. Let $A \in M_{n,k}(\mathbb{R})$ have rank k . How many ways can we write $A = QR$ where $Q \in M_{n,k}(\mathbb{R})$ has orthonormal columns and $R \in M_k(\mathbb{R})$ is upper-triangular? (We do not assume R has positive diagonal entries here.)
16. Use Gram–Schmidt orthonormalization to prove that any $A \in M_{n,k}(\mathbb{C})$ can be factored as $A = QR$, where $Q \in M_n(\mathbb{C})$ is unitary and $R \in M_{n,k}(\mathbb{C})$ is upper-triangular.
17. Formulate and prove a uniqueness result for the factorization $A = QR$ of an arbitrary matrix $A \in M_{n,k}(\mathbb{C})$ described in §9.3.
18. *Modified Gram–Schmidt Algorithm.* Given a linearly independent list (v_1, \dots, v_k) in an inner product space V , the Gram–Schmidt orthonormalization algorithm in §9.2 produces an orthonormal list (u_1, \dots, u_k) . Show that the following algorithm will also transform (v_1, \dots, v_k) into (u_1, \dots, u_k) . Loop from $j = 1$ to k . Replace the current vector v_j by $v_j / \|v_j\|$, and then for each r from $j+1$ to k , replace v_r by $v_r - \langle v_r, v_j \rangle v_j$. Output the final values of (v_1, \dots, v_k) . (This

version of Gram–Schmidt can be shown to be more stable numerically compared to the original version.)

19. Estimate the number of real multiplications, divisions, and square root extractions needed to compute the QR factorization of a matrix $A \in M_n(\mathbb{R})$ via the Gram–Schmidt orthonormalization algorithm.
20. Estimate the number of real multiplications, divisions, and square root extractions needed to compute the QR factorization of a matrix $A \in M_n(\mathbb{R})$ via Householder’s algorithm.
21. Describe how to adapt the Gram–Schmidt algorithm to transform a linearly independent list (v_1, \dots, v_k) into an *orthogonal* list (u_1, \dots, u_k) using a computation that avoids extractions of square roots.
22. (a) Draw a picture in \mathbb{R}^2 illustrating the action of $T_{(-4,3)}$ on each of these vectors: $(-4, 3); (3, 4); (1, 0); (0, 1); (1, -2)$. (b) Confirm your answers to (a) algebraically by multiplying each column vector by $Q_{(-4,3)}$.
23. Find a Householder matrix $Q_v \in M_3(\mathbb{R})$ that fixes each point in the plane $2x_1 - 3x_2 + x_3 = 0$. Draw a sketch showing this plane and the vectors v , $w = (1, -1, 1)$, and $Q_v w$.
24. (a) Explain why the lemma in §9.4 fails to hold without the hypothesis $\|x\| = \|y\|$. (b) Does the lemma hold without the hypothesis $\langle x, y \rangle \in \mathbb{R}$? (c) Suppose $x = y$. Describe all $v \in \mathbb{C}^n$ with $Q_v x = y$.
25. For each $x, y \in \mathbb{C}^n$, find v and Q_v such that $Q_v x = y$ or explain why this cannot be done. (a) $x = (5, 12)$, $y = (13, 0)$. (b) $x = (1, i)$, $y = (\sqrt{2}, 0)$. (c) $x = (i, 1)$, $y = (\sqrt{2}, 0)$. (d) $x = (3, 2, 6)$, $y = (2, 6, 3)$. (e) $x = (3, 2, 3)$, $y = (2, 3, 2)$. (f) $x = (3, 2, 6)$, $y = (b, 0, 0)$ for different choices of $b \in \mathbb{C}$.
26. For each matrix in Exercise 14, use Householder’s algorithm to compute the QR factorization of A . In each case, write Q as an explicit product of Householder reflections and a diagonal unitary matrix.
27. Apply Gram–Schmidt orthonormalization to the columns of the matrix A in §9.5 to obtain the matrices Q and R in (9.2).
28. *Givens Rotations.* For any $i \neq j$ in $[n]$ and $\theta \in [0, 2\pi)$, define the *Givens rotation matrix* $G = G(i, j; \theta) \in M_n(\mathbb{R})$ by letting $G(i, i) = G(j, j) = \cos \theta$, $G(j, i) = \sin \theta$, $G(i, j) = -\sin \theta$, $G(k, k) = 1$ for all $k \neq i, j$, and letting all other entries of G be zero. Left-multiplication by $G(i, j; \theta)$ rotates the plane spanned by e_i and e_j counterclockwise through an angle of θ . (a) Show that $G(i, j; \theta)$ is orthogonal with inverse $G(i, j; -\theta)$. (b) Given a matrix $A \in M_n(\mathbb{R})$ with $A(j, i) \neq 0$, find a specific θ such that $G(i, j; \theta)A$ has j , i -entry zero.
29. *Givens’ QR Algorithm.* Describe an algorithm for computing a factorization $A = QR$ of a matrix $A \in M_{n,k}(\mathbb{R})$ that reduces A to an upper-triangular matrix R by left-multiplying by a sequence of Givens rotation matrices (see Exercise 28). How many matrices are needed in general?
30. (a) Show that every unitary matrix $Q \in M_n(\mathbb{C})$ can be factored as $Q = Q_{v_1} \cdots Q_{v_{n-1}} D$, where D is a diagonal unitary matrix and each Q_{v_j} is a Householder matrix. (b) Can the matrix D in (a) be omitted? Explain.
31. In numerical computations, one source of inaccuracy is *subtractive cancellation*, in which two nearly equal real numbers are subtracted. Explain how the minus sign in the formula for y in §9.5 avoids subtractive cancellation in the computation of the first component of v_1 .

32. Find Doolittle's LU factorization of each matrix using the recursive formulas in §9.6. (a) $\begin{bmatrix} 3 & -2 \\ -12 & 13 \end{bmatrix}$; (b) $\begin{bmatrix} 2 & -3 & 1 \\ -4 & 10 & 6 \\ 6 & -13 & -3 \end{bmatrix}$; (c) $\begin{bmatrix} 2 & 3 & -2 & 4 \\ 1 & 6.5 & -9 & 3 \\ -5 & 2.5 & -9 & -8 \\ -2 & -3 & 5 & -4 \end{bmatrix}$.
33. Find an LU factorization of each matrix using Gaussian elimination. (a) $\begin{bmatrix} 2 & 8 \\ 6 & 19 \end{bmatrix}$; (b) $\begin{bmatrix} 3 & 2 & 1 \\ 2 & 10/3 & 5/3 \\ -5 & 8/3 & 7/3 \end{bmatrix}$; (c) $\begin{bmatrix} 2 & 10 & -13 & 8 \\ -10 & -35 & 58 & -46 \\ -6 & -24 & 39.2 & -35.4 \\ 0.8 & -41 & 16.4 & 21.4 \end{bmatrix}$.
34. Justify the existence and uniqueness assertions for the variants of the LU factorization described in the last paragraph of §9.6.
35. State and prove a theorem regarding factorizations of square matrices of the form $A = UL$, where U is upper-triangular and L is lower-triangular. Give one proof by modifying the formulas in §9.6, and a second proof that deduces UL factorizations from LU factorizations.
36. Show that for any field F and any permutation matrix $P \in M_n(F)$, $P^T = P^{-1}$.
37. For each matrix A , find a factorization $PA = LU$ or $PAQ = LU$ as in §9.9. (a) $A = \begin{bmatrix} 0 & 2 \\ 3 & 1 \end{bmatrix}$; (b) $A = \begin{bmatrix} 0 & 0 & 2 \\ 1 & -1 & 2 \\ 3 & 3 & 3 \end{bmatrix}$; (c) $A = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 3 \\ 0 & 5 & 0 & 0 \end{bmatrix}$.
38. Extend the results on LU factorizations and permuted LU factorizations to rectangular matrices $A \in M_{m,n}(F)$.
39. Without using Gaussian elimination or permuted LU factorizations, prove by induction on n that for every invertible $A \in M_n(F)$, there exists a permutation matrix $P \in M_n(F)$ such that $\det((PA)[k]) \neq 0$ for $1 \leq k \leq n$.
40. Suppose $A \in M_n(\mathbb{C})$ has rank r and $\det(A[k]) \neq 0$ for $1 \leq k \leq r$. Prove that A has an LU factorization in which L is lower-unitriangular and U is upper-triangular with the last $n - r$ diagonal entries equal to zero.
41. Show that any $A \in M_n(F)$ can be factored as $A = LQU$, where L is lower-triangular and invertible, U is upper-triangular and invertible, and $Q \in M_n(F)$ has at most one 1 in each row and column, and all other entries of Q are zero. (Use elementary matrices to reduce A .)
42. Find a Cholesky factorization of each positive semidefinite matrix. (a) $\begin{bmatrix} 4 & -2i \\ 2i & 2 \end{bmatrix}$; (b) $\begin{bmatrix} 9 & 15 & 0 \\ 15 & 26 & 1 \\ 0 & 1 & 5 \end{bmatrix}$; (c) $\begin{bmatrix} 16 & -4 & -12 & 4 \\ -4 & 1 & 3 & -1 \\ -12 & 3 & 19 & -1 \\ 4 & -1 & -1 & 4 \end{bmatrix}$.
43. Give an example to show that the Cholesky factorization of a non-invertible positive semidefinite matrix need not be unique.
44. (a) Count the number of multiplications, divisions, and square root extractions in \mathbb{C} needed to compute the Cholesky factorization of an $n \times n$ positive definite matrix using the formulas in §9.10. (b) How do the operation counts in (a) compare to those needed to find the LU decomposition of an $n \times n$ matrix using the formulas in §9.6?

45. Consider the (inconsistent) linear system $x + y = 3$, $2x - y = 0$, $x - 2y = 0$.
 (a) Use the normal equations to find the least squares approximate solution to this system. (b) Graph the three equations in the system and the least squares solution in one picture.
46. Consider the linear system

$$\left\{ \begin{array}{rcl} 3x & +7y & -2z = 4 \\ 5x & -3y & +2z = 8 \\ -x & & -5z = 2 \\ -2x & +y & +3z = -1 \\ 7x & -4y & +z = 0 \end{array} \right.$$

Find the least squares approximate solution to the system by setting up and solving the normal equations.

47. (a) Given the system of two equations $ax = b$ and $cx = d$ with $(a, c) \neq (0, 0)$, what is the least squares approximate solution to this system? (b) Given n equations $a_j x = b_j$ with $a_1, \dots, a_n, b_1, \dots, b_n \in \mathbb{C}$ and not all a_j are zero, what is the least squares approximate solution to this system?
48. Suppose we want the least squares approximate solution to $Ax = b$, where $A \in M_{n,k}(\mathbb{C})$ has rank k , and we know a factorization $A = QR$ with the properties stated in §9.3. Show that this solution can be found efficiently by solving $Rx = Q^*b$.
49. Given n real data points $(x_1, y_1), \dots, (x_n, y_n)$, suppose we want to find real parameters m and b such that the line $y = mx + b$ “best approximates” the given data. Formulate this as a least squares problem, and solve the normal equations to obtain specific formulas for m and b . Specify exactly what quantity is being minimized by the optimal choice of m and b .
50. Follow the proof in §9.12 to find a singular value decomposition for each matrix.

$$(a) \begin{bmatrix} 1.824 & 1.032 \\ -1.968 & 2.176 \end{bmatrix}; \quad (b) \begin{bmatrix} 3 & 3 & 3 \\ -1 & -1 & -1 \end{bmatrix}; \quad (c) \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \end{bmatrix};$$

$$(d) \begin{bmatrix} 2 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 \\ 2 & 0 & 2 & 0 \\ 0 & 0 & 0 & i \end{bmatrix}.$$

51. *Penrose Properties.* (a) Prove that for every $A \in M_{m,n}(\mathbb{C})$, there exists at most one matrix $B \in M_{n,m}(\mathbb{C})$ with these four properties, called the *Penrose properties*.

$$\text{I: } ABA = A; \quad \text{II: } BAB = B; \quad \text{III: } (AB)^* = AB; \quad \text{IV: } (BA)^* = BA.$$

[If B and C both have these properties, show $B = (BA)^*(CA)^*C(AC)^*(AB)^* = C$.] (b) Show that if $A \in M_n(\mathbb{C})$ is invertible, then $B = A^{-1}$ has the four Penrose properties.

52. *Pseudoinverses.* Suppose $A \in M_{m,n}(\mathbb{C})$ has singular value decomposition $A = VDU^*$. Let $D' \in M_{n,m}(\mathbb{C})$ be the matrix obtained from D^T by inverting the nonzero diagonal entries. The matrix $A^+ = UD'V^* \in M_{n,m}(\mathbb{C})$ is called the *pseudoinverse of A* . (a) Show that the pseudoinverse does not depend on the choice of singular value decomposition of A , by showing that A^+ has the four Penrose properties in Exercise 51. (b) Without using the Penrose properties, show that $A^+ = A^{-1}$ when $A \in M_n(\mathbb{C})$ is invertible.

53. For each matrix A in Exercise 50, find the pseudoinverse A^+ (see Exercise 52). Compute AA^+ and A^+A in each case.
54. Let $A = VDU^*$ be a singular value decomposition of $A \in M_{m,n}(\mathbb{C})$ such that D has k nonzero entries on the diagonal. (a) Show that the last $n - k$ columns of U are an orthonormal basis of the null space of A . (b) Show that the first k columns of V are an orthonormal basis of the image of A .
55. Let $A \in M_{m,n}(\mathbb{C})$. Use singular value decompositions to prove that $\sup\{\|Ax\| : x \in \mathbb{C}^n, \|x\| = 1\}$ is the largest singular value of A . (Here, $\|x\| = \|x\|_2 = \sqrt{\langle x, x \rangle}$.)

Iterative Algorithms in Numerical Linear Algebra

In applications of linear algebra to science, engineering, and other areas, one often needs to find numerical solutions to a huge linear system of equations or to approximate the eigenvalues of a large square matrix. The branch of numerical analysis called *numerical linear algebra* studies algorithms for solving such problems efficiently while minimizing the effects of rounding errors.

Consider the problem of solving the linear system $Ax = b$, where $A \in M_n(\mathbb{C})$ is a given invertible matrix, $b \in \mathbb{C}^n$ is a given column vector, and $x \in \mathbb{C}^n$ is a column vector of unknowns. Gaussian elimination and related algorithms solve this system by a finite sequence of steps, which (in the absence of roundoff errors) will ultimately terminate with the exact solution. This chapter discusses a different class of algorithms for solving such problems, called *iterative algorithms*. The idea here is to start with some initial approximation x_0 to the sought-for solution x , and then to apply some recursive formula $x_{k+1} = f(x_k)$ (depending on the problem instance) to compute an infinite sequence $x_0, x_1, x_2, \dots, x_k, \dots$ of approximate solutions. We can then ask if this sequence converges (under an appropriate definition of convergence) to the true solution x , and try to find bounds on “how far away” a given approximation x_k is to x .

In order to give a precise mathematical analysis of such algorithms, we need a formal way to measure the *distance* between two column vectors in \mathbb{C}^n . We will define this distance using the idea of a *vector norm*. The related concept of a *matrix norm* will help us understand how multiplication by a fixed matrix affects the distances between vectors.

The analysis of iterative algorithms constitutes a vast subfield of numerical linear algebra, and we only have space to give a brief taste of this subject in this chapter. We discuss three basic iterative methods for solving linear systems (the algorithms of Richardson, Jacobi, and Gauss–Seidel) and analyze them using norms and the notion of the spectral radius of a matrix. We also study the power method, which can be used to approximate the largest eigenvalue of a given square matrix. To find other eigenvalues of the matrix, one can use variations of the power method or a technique called *deflation*.

10.1 Richardson’s Algorithm

As a first illustration of an iterative method for solving the linear system $Ax = b$, we describe *Richardson’s algorithm*. The input to the algorithm is a matrix $A \in M_n(\mathbb{C})$, a column vector $b \in \mathbb{C}^n$, and a vector $x_0 \in \mathbb{C}^n$ that represents an initial guess or approximation to the true solution vector x . (We may take $x_0 = 0$ if we have no initial information about x .) The algorithm proceeds by forming a sequence of approximate solutions $x_0, x_1, x_2, \dots, x_k, \dots$, where for all $k \geq 0$,

$$x_{k+1} = x_k + (b - Ax_k).$$

Intuitively, having already computed x_k , we find the *output error vector* $b - Ax_k$ and add this vector to x_k to obtain the next input vector x_{k+1} . Ideally, we would reach the true

solution x from x_k by adding the *input error vector* $x - x_k$ to x_k , but we do not know this error vector. However, if the matrix A is “close” to the identity matrix in some sense, then we would expect the output error vector $b - Ax_k = A(x - x_k)$ to be “close” to the input error vector. In this case, we would expect the approximate solutions x_k to approach the true solution x . We make this informal analysis rigorous in §10.11, where we also derive a precise bound for the size of the error vector $x - x_k$.

Here is a small example to illustrate the execution of Richardson’s algorithm. Let $A = \begin{bmatrix} 1.1 & 0.3 \\ -0.4 & 0.8 \end{bmatrix}$ and $b = \begin{bmatrix} 2 \\ -3 \end{bmatrix}$. Starting with $x_0 = 0$, we successively compute

$$\begin{aligned} x_1 &= x_0 + (b - Ax_0) = \begin{bmatrix} 2 \\ -3 \end{bmatrix}, & x_2 &= x_1 + (b - Ax_1) = \begin{bmatrix} 2.7 \\ -2.8 \end{bmatrix}, \\ x_3 &= x_2 + (b - Ax_2) = \begin{bmatrix} 2.57 \\ -2.48 \end{bmatrix}, & x_4 &= x_3 + (b - Ax_3) = \begin{bmatrix} 2.487 \\ -2.468 \end{bmatrix}, \\ x_5 &= x_4 + (b - Ax_4) = \begin{bmatrix} 2.4917 \\ -2.4988 \end{bmatrix}, & x_6 &= x_5 + (b - Ax_5) = \begin{bmatrix} 2.50047 \\ -2.50308 \end{bmatrix}. \end{aligned}$$

On the other hand, by Gaussian elimination or Cramer’s rule, we find the exact solution is $x = \begin{bmatrix} 2.5 \\ -2.5 \end{bmatrix}$. We see that the sequence is converging rapidly to the exact solution.

However, suppose we take $A = \begin{bmatrix} 2 & 3 \\ 1 & 4 \end{bmatrix}$ and $b = \begin{bmatrix} 2 \\ -3 \end{bmatrix}$. Starting the algorithm at $x_0 = 0$, we compute the sequence of approximations

$$\left(\begin{bmatrix} 2 \\ -3 \end{bmatrix}, \begin{bmatrix} 9 \\ 4 \end{bmatrix}, \begin{bmatrix} -19 \\ -24 \end{bmatrix}, \begin{bmatrix} 93 \\ 88 \end{bmatrix}, \begin{bmatrix} -355 \\ -360 \end{bmatrix}, \begin{bmatrix} 1437 \\ 1432 \end{bmatrix}, \dots \right),$$

which does not appear to be converging. (The true solution here is $\begin{bmatrix} 3.4 \\ -1.6 \end{bmatrix}$.)

10.2 Jacobi’s Algorithm

We introduce the next iterative algorithm for solving $Ax = b$, called *Jacobi’s algorithm*,

with a specific example. Suppose $A = \begin{bmatrix} 5 & 1 & -1 \\ -1 & 5 & 2 \\ 2 & 2 & 4 \end{bmatrix}$, $b = \begin{bmatrix} 4 \\ -1 \\ 3 \end{bmatrix}$, and $x = \begin{bmatrix} r \\ s \\ t \end{bmatrix}$. We

are trying to solve a system of three linear equations in three unknowns. The basic idea is to solve the i ’th equation for the i ’th unknown, as shown here:

$$\left\{ \begin{array}{lcl} 5r + s - t & = & 4 \\ -r + 5s + 2t & = & -1 \\ 2r + 2s + 4t & = & 3 \end{array} \right. \Rightarrow \left\{ \begin{array}{lcl} r & = & (4 - s + t)/5 \\ s & = & (-1 + r - 2t)/5 \\ t & = & (3 - 2r - 2s)/4. \end{array} \right.$$

We use the rewritten equations to define a sequence of vectors $x_k = [r_k \ s_k \ t_k]^T$ by choosing any initial vector $x_0 = [r_0 \ s_0 \ t_0]^T$, and then computing

$$\left\{ \begin{array}{lcl} r_{k+1} & = & (4 - s_k + t_k)/5 \\ s_{k+1} & = & (-1 + r_k - 2t_k)/5 \\ t_{k+1} & = & (3 - 2r_k - 2s_k)/4 \end{array} \right. \quad \text{for all } k \geq 0. \quad (10.1)$$

If we start with $x_0 = 0$, we compute the following sequence of approximate solutions:

$$\begin{bmatrix} 0.8 \\ -0.2 \\ 0.75 \end{bmatrix}, \begin{bmatrix} 0.99 \\ -0.34 \\ 0.45 \end{bmatrix}, \begin{bmatrix} 0.958 \\ -0.182 \\ 0.425 \end{bmatrix}, \begin{bmatrix} 0.921 \\ -0.178 \\ 0.362 \end{bmatrix}, \begin{bmatrix} 0.908 \\ -0.161 \\ 0.379 \end{bmatrix}, \begin{bmatrix} 0.908 \\ -0.170 \\ 0.376 \end{bmatrix}, \begin{bmatrix} 0.909 \\ -0.169 \\ 0.381 \end{bmatrix}, \dots$$

These vectors appear to be converging rapidly to the exact solution $x = [0.91 \ -0.17 \ 0.38]^T$.

Now we describe how Jacobi's algorithm works in general. The input to the algorithm is a matrix $A \in M_n(\mathbb{C})$, a column vector $b \in \mathbb{C}^n$, and an initial approximation $x_0 \in \mathbb{C}^n$. This algorithm requires that all diagonal entries $A(i, i)$ be nonzero; in fact, it is preferable for the diagonal entry in each row to be large in magnitude compared to the other entries in its row. If the initial matrix does not satisfy this requirement, we may adjust A by permuting rows and columns to make the requirement hold. Row permutations correspond to reordering the linear equations in the system $Ax = b$, whereas column permutations correspond to reordering the unknown components of x .

Write $A = Q + R$, where Q consists of the diagonal entries of A , and R consists of the off-diagonal entries of A . Formally, let $Q(i, i) = A(i, i)$ for $i \in [n]$, let all other entries of Q be zero, let $R(i, j) = A(i, j)$ for $i, j \in [n]$ with $i \neq j$, and let all other entries of R be zero. Since Q is a diagonal matrix with nonzero entries on the diagonal, we can quickly compute Q^{-1} by inverting each diagonal entry. We now compute the sequence of approximations $x_0, x_1, x_2, \dots, x_k, \dots$ by setting

$$x_{k+1} = Q^{-1}(b - Rx_k) \quad \text{for all } k \geq 0. \quad (10.2)$$

To see that this matches the description of the algorithm given in the example above, note that the i 'th row of $Ax = b$ is the equation

$$b(i) = \sum_{j=1}^n A(i, j)x(j) = A(i, i)x(i) + \sum_{j \neq i} A(i, j)x(j) = Q(i, i)x(i) + \sum_{j=1}^n R(i, j)x(j).$$

(We use parentheses to indicate components of a vector to avoid confusion with the subscripts in the sequence of approximations x_0, x_1, x_2, \dots) Solving the i 'th equation for $x(i)$ and using the resulting formula as the update rule for obtaining $x_{k+1}(i)$ from x_k , we get

$$x_{k+1}(i) = Q(i, i)^{-1} \left[b(i) - \sum_{j=1}^n R(i, j)x_k(j) \right]$$

for all $i \in [n]$. These formulas are equivalent to the matrix equation (10.2).

We analyze the convergence properties of Jacobi's algorithm in §10.12.

10.3 Gauss-Seidel Algorithm

The *Gauss-Seidel algorithm* is a variation of Jacobi's algorithm that uses the updated value of each unknown variable as soon as that value is computed. We illustrate the idea with the same 3×3 system $Ax = b$ considered in §10.2. The key modification to the previous example is that the update rules in (10.1) are replaced by

$$\begin{cases} r_{k+1} &= (4 & -s_k & +t_k)/5 \\ s_{k+1} &= (-1 & +r_{k+1} & -2t_k)/5 \\ t_{k+1} &= (3 & -2r_{k+1} & -2s_{k+1})/4 \end{cases} \quad \text{for all } k \geq 0. \quad (10.3)$$

The update rule for r_{k+1} is the same as before. When we compute s_{k+1} , we have already computed the new value r_{k+1} but we have not yet computed t_{k+1} . So we use r_{k+1} and t_k on the right side of the update rule for s_{k+1} . When we compute t_{k+1} , we use the newly computed values r_{k+1} and s_{k+1} on the right side instead of the old values r_k and s_k .

Using the Gauss–Seidel formulas and starting with $x_0 = 0$, we compute the following sequence of approximate solutions:

$$\begin{bmatrix} 0.8 \\ -0.04 \\ 0.37 \end{bmatrix}, \begin{bmatrix} 0.882 \\ -0.172 \\ 0.395 \end{bmatrix}, \begin{bmatrix} 0.913 \\ -0.175 \\ 0.381 \end{bmatrix}, \begin{bmatrix} 0.911 \\ -0.170 \\ 0.379 \end{bmatrix}, \begin{bmatrix} 0.910 \\ -0.170 \\ 0.380 \end{bmatrix}, \dots$$

This sequence appears to be converging to the exact solution $x = [0.91 \ -0.17 \ 0.38]^T$ even faster than the sequence produced by the Jacobi method.

To solve a general system $Ax = b$ using the Gauss–Seidel algorithm, we require that every diagonal entry of A be nonzero (which may require preprocessing A by permuting rows and/or columns). Starting with any initial vector x_0 , we compute x_1, x_2, \dots , via the update rules

$$x_{k+1}(i) = A(i, i)^{-1} \left[b(i) - \sum_{j < i} A(i, j)x_{k+1}(j) - \sum_{j > i} A(i, j)x_k(j) \right] \quad \text{for all } k \geq 0, \quad (10.4)$$

which are computed in order as i increases from 1 to n .

Next we want to give a matrix formulation of the Gauss–Seidel iterative procedure. To see how this is done, let us first rewrite (10.3) by moving all terms with subscript $k + 1$ to the left side. We obtain

$$\begin{bmatrix} 5 & 0 & 0 \\ -1 & 5 & 0 \\ 2 & 2 & 4 \end{bmatrix} \begin{bmatrix} r_{k+1} \\ s_{k+1} \\ t_{k+1} \end{bmatrix} = \begin{bmatrix} 4 \\ -1 \\ 3 \end{bmatrix} + \begin{bmatrix} 0 & -1 & 1 \\ 0 & 0 & -2 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} r_k \\ s_k \\ t_k \end{bmatrix}.$$

The matrix on the left consists of the entries in the original matrix A lying on or below the main diagonal, whereas the matrix on the right is found by negating the entries of A above the main diagonal. For the general system $Ax = b$, define Q and R in $M_n(\mathbb{C})$ by setting $Q(i, j) = A(i, j)$ for all $i \geq j$ in $[n]$, $Q(i, j) = 0$ for $i < j$, and $R = A - Q$. Multiplying (10.4) by $A(i, i)$ and adding $\sum_{j < i} A(i, j)x_{k+1}(j)$ to both sides, we see that the equations (10.4) for $i \in [n]$ are equivalent to the matrix update rule

$$Qx_{k+1} = b - Rx_k \quad \text{for all } k \geq 0.$$

When using this update rule, it should be kept in mind that we do not find x_{k+1} by explicitly computing Q^{-1} and applying this matrix to $b - Rx_k$. Rather, we compute $b - Rx_k$ and then “backsolve” for x_{k+1} as indicated in (10.4). Note that (10.2) can also be written in the form $Qx_{k+1} = b - Rx_k$, but with a different choice of Q and R . In both cases, a key property of Q is that $Qx_{k+1} = y$ can be quickly solved for x_{k+1} in much less time than the n^3 steps required for a general matrix inversion.

We study the convergence properties of the Gauss–Seidel algorithm in §10.14.

10.4 Vector Norms

To analyze the convergence of iterative algorithms, we need the idea of the distance between two vectors in a vector space. To define this, we first introduce the concept of the length

(or norm) of a vector. Given a real or complex vector space V , a *norm* on V is a function $N : V \rightarrow \mathbb{R}$, denoted $N(x) = \|x\|$ for $x \in V$, which must satisfy the following axioms. First, $0 \leq \|x\| < \infty$ for all $x \in V$, and $\|x\| = 0$ iff $x = 0$. Second, for all $x \in V$ and all scalars c , $\|cx\| = |c| \cdot \|x\|$. Third, for all $x, y \in V$, the *triangle inequality* $\|x + y\| \leq \|x\| + \|y\|$ must hold. A vector space V together with a specific norm N on V is called a *normed vector space*.

For example, let V be the real vector space \mathbb{R}^n . For $x = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$, define the *sup norm* by setting

$$\|x\|_\infty = \max(|x_1|, |x_2|, \dots, |x_n|). \quad (10.5)$$

Let us check the axioms for a normed vector space. Fix $x, y \in \mathbb{R}^n$ and $c \in \mathbb{R}$. First, since $|x_i| \geq 0$ for all i , the maximum of these n numbers is also nonnegative and finite; and the maximum will equal zero iff all $x_i = 0$ iff $x = 0$. So $\|x\|_\infty \geq 0$ with equality iff $x = 0$. Second, $|cx_i| = |c| \cdot |x_i|$ for each i , so

$$\|cx\|_\infty = \max(|c| \cdot |x_1|, \dots, |c| \cdot |x_n|) = |c| \max(|x_1|, \dots, |x_n|) = |c| \cdot \|x\|_\infty.$$

Third, the triangle inequality in \mathbb{R} tells us that $|x_i + y_i| \leq |x_i| + |y_i|$ for $1 \leq i \leq n$. Now, $|x_i| \leq \max_{1 \leq j \leq n} |x_j| = \|x\|_\infty$ and $|y_i| \leq \max_{1 \leq j \leq n} |y_j| = \|y\|_\infty$, so

$$|x_i + y_i| \leq \|x\|_\infty + \|y\|_\infty.$$

This inequality holds for $1 \leq i \leq n$, and therefore

$$\|x + y\|_\infty = \max(|x_1 + y_1|, \dots, |x_n + y_n|) \leq \|x\|_\infty + \|y\|_\infty,$$

as needed. An identical proof shows that \mathbb{C}^n is a complex normed vector space using the sup norm; in this case, $|x_j|$ denotes the magnitude of the complex number x_j .

Now we describe another norm on the real vector space \mathbb{R}^n . This is the *1-norm*, defined by

$$\|x\|_1 = |x_1| + |x_2| + \cdots + |x_n| \quad (10.6)$$

for all $x = (x_1, \dots, x_n) \in \mathbb{R}^n$. The first two axioms are readily verified; let us check the triangle inequality. Given $x, y \in \mathbb{R}^n$, we have $|x_i + y_i| \leq |x_i| + |y_i|$ in \mathbb{R} for $1 \leq i \leq n$. Adding these n inequalities, we obtain

$$\|x + y\|_1 = \sum_{i=1}^n |x_i + y_i| \leq \sum_{i=1}^n |x_i| + \sum_{i=1}^n |y_i| = \|x\|_1 + \|y\|_1.$$

A third example of a norm on \mathbb{R}^n is the *Euclidean norm* or *2-norm*, defined by

$$\|x\|_2 = \sqrt{|x_1|^2 + |x_2|^2 + \cdots + |x_n|^2} = \sqrt{x^* x} \quad (10.7)$$

for $x = (x_1, \dots, x_n) \in \mathbb{R}^n$. The reader may check that the norm axioms are satisfied; the proof of the triangle inequality follows from the Cauchy–Schwarz inequality (see Exercise 26). Similarly, one can verify that formulas (10.6) and (10.7) define norms on the complex vector space \mathbb{C}^n .

A vector u in a normed vector space V is called a *unit vector* iff $\|u\| = 1$. Given any nonzero $v \in V$, we claim $\|v\|^{-1}v$ is a unit vector. For, setting $c = \|v\|^{-1}$, we have $\|cv\| = |c| \cdot \|v\| = \|v\|^{-1}\|v\| = 1$.

10.5 Metric Spaces

A *metric space* is a set X together with a *distance function* (or *metric*) $d : X \times X \rightarrow \mathbb{R}$, which must satisfy the following conditions. First, for all $x, y \in X$, $0 \leq d(x, y) < \infty$, and $d(x, y) = 0$ iff $x = y$. Second, for all $x, y \in X$, $d(x, y) = d(y, x)$. Third, for all $x, y, z \in X$, the *triangle inequality* $d(x, z) \leq d(x, y) + d(y, z)$ holds. We think of the nonnegative real number $d(x, y)$ as the *distance* from x to y .

Any normed vector space V becomes a metric space if we define $d(x, y) = \|x - y\|$ for $x, y \in V$. To check this, fix $x, y, z \in V$. We have $0 \leq \|x - y\| < \infty$, with $\|x - y\| = 0$ iff $x - y = 0$ iff $x = y$, so the first axiom for a metric space holds. We have

$$d(y, x) = \|y - x\| = \|(-1)(x - y)\| = |-1| \cdot \|x - y\| = \|x - y\| = d(x, y),$$

so the second axiom holds. Finally,

$$d(x, z) = \|x - z\| = \|(x - y) + (y - z)\| \leq \|x - y\| + \|y - z\| = d(x, y) + d(y, z),$$

so the third axiom holds. A metric arising from a norm in this way has two additional properties related to the vector space structure of V . First, for all $x, y, z \in V$, $d(x + z, y + z) = d(x, y)$; this holds since $d(x + z, y + z) = \|(x + z) - (y + z)\| = \|x - y\| = d(x, y)$. Second, for all $x, y \in V$ and all scalars c , $d(cx, cy) = |c|d(x, y)$; this holds since $d(cx, cy) = \|cx - cy\| = \|c(x - y)\| = |c| \cdot \|x - y\| = |c|d(x, y)$. We summarize these properties by saying that the metric is *translation-invariant* and *respects dilations*.

10.6 Convergence of Sequences

In any metric space (X, d) , we can define the idea of a convergent sequence. Given $x \in X$ and given a sequence $(x_0, x_1, x_2, \dots) = (x_k : k \geq 0)$ of points in X , we say that the sequence (x_k) *converges* to x iff for every $\epsilon > 0$, there exists $k_0 \in \mathbb{N}$ such that for all $k \geq k_0$, $d(x_k, x) < \epsilon$. Informally, this definition says that the sequence gets arbitrarily close to x (and stays close) for terms far enough out in the sequence. When this condition holds, we write $\lim_{k \rightarrow \infty} x_k = x$, and we also write $x_k \rightarrow x$ as $k \rightarrow \infty$.

Not all sequences in a metric space (X, d) converge. But if a sequence $(x_k : k \geq 0)$ does converge to some limit x , then x is unique. For suppose $x_k \rightarrow x$ and also $x_k \rightarrow y$ for some $x, y \in X$. Given $\epsilon > 0$, choose $k_1, k_2 \in \mathbb{N}$ so that $d(x_k, x) < \epsilon/2$ for all $k \geq k_1$ and $d(x_k, y) < \epsilon/2$ for all $k \geq k_2$. Then taking $k = \max(k_1, k_2)$, $d(x, y) \leq d(x, x_k) + d(x_k, y) < \epsilon$. Since this holds for every positive ϵ , we must have $d(x, y) = 0$, so that $x = y$ as needed.

Now suppose V is a normed vector space; we always assume V has the metric $d(x, y) = \|x - y\|$ defined using the norm. Given two sequences $(v_k : k \geq 0)$ and $(w_k : k \geq 0)$ in V , suppose $v_k \rightarrow v$ and $w_k \rightarrow w$ for some $v, w \in V$. We assert that $v_k + w_k \rightarrow v + w$. To prove this, fix $\epsilon > 0$. Choose $k_1, k_2 \in \mathbb{N}$ such that $d(v_k, v) < \epsilon/2$ for all $k \geq k_1$ and $d(w_k, w) < \epsilon/2$ for all $k \geq k_2$. Then for all $k \geq \max(k_1, k_2)$, the triangle inequality gives

$$\begin{aligned} d(v_k + w_k, v + w) &= \|(v_k + w_k) - (v + w)\| = \|(v_k - v) + (w_k - w)\| \\ &\leq \|v_k - v\| + \|w_k - w\| < \epsilon/2 + \epsilon/2 = \epsilon. \end{aligned}$$

So $\lim_{k \rightarrow \infty} (v_k + w_k) = v + w$.

Next, assume c is a scalar and $v_k \rightarrow v$; we assert that $cv_k \rightarrow cv$. This result is evident

when $c = 0$, so assume $c \neq 0$. Given $\epsilon > 0$, choose k_0 so that $d(v_k, v) < \epsilon/|c|$ for all $k \geq k_0$. Then for all $k \geq k_0$, we compute $d(cv_k, cv) = |c|d(v_k, v) < \epsilon$. So $\lim_{k \rightarrow \infty} (cv_k) = cv$. More generally, one can show that if $c_k \rightarrow c$ in \mathbb{R} or \mathbb{C} and $v_k \rightarrow v$ in V , then $c_k v_k \rightarrow cv$ in V (Exercise 50). One can also show by induction that if $(v_k^{(1)} : k \geq 0), \dots, (v_k^{(m)} : k \geq 0)$ are m sequences converging respectively to $v^{(1)}, \dots, v^{(m)}$, and if c_1, \dots, c_m are fixed scalars, then

$$\lim_{k \rightarrow \infty} (c_1 v_k^{(1)} + \dots + c_m v_k^{(m)}) = c_1 v^{(1)} + \dots + c_m v^{(m)}.$$

10.7 Comparable Norms

As we have seen, a given vector space V may possess many different norms and hence many different metrics. The definition of a convergent sequence depends on the metric used, so it is conceivable that a given sequence $(v_k : k \geq 0)$ might converge to some v relative to one norm, but not converge to v relative to some other norm. However, when V is finite-dimensional, we will show that this does not happen. In other words, for any two norms $\|\cdot\|$ and $\|\cdot\|'$ on V , we will see that $v_k \rightarrow v$ relative to $\|\cdot\|$ iff $v_k \rightarrow v$ relative to $\|\cdot\|'$. So for studying questions of convergence, we can use whichever norm is most convenient for calculations.

Let us say that two norms $\|\cdot\|$ and $\|\cdot\|'$ on a vector space V are *comparable* iff there exist real constants C, D with $0 < C, D < \infty$ such that for all $x \in V$, $\|x\| \leq C\|x\|'$ and $\|x\|' \leq D\|x\|$. For example, let us show that the sup norm and the 2-norm are comparable on $V = \mathbb{R}^n$. It suffices to prove that for all $x = (x_1, \dots, x_n) \in \mathbb{R}^n$,

$$\|x\|_\infty \leq \|x\|_2 \leq \sqrt{n}\|x\|_\infty$$

(here $C = 1$ and $D = \sqrt{n}$). For the first inequality, observe that $|x_i| = \sqrt{|x_i|^2} \leq \sqrt{\sum_{j=1}^n |x_j|^2} = \|x\|_2$ for all i . Taking the maximum over all i gives $\|x\|_\infty \leq \|x\|_2$. Next, fix an index k such that $|x_k| = \max_i |x_i|$. Then $|x_i|^2 \leq |x_k|^2$ for all i , so

$$\|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2} \leq \sqrt{\sum_{i=1}^n |x_k|^2} = \sqrt{n}|x_k| = \sqrt{n}\|x\|_\infty.$$

Similarly, the reader may prove that the sup norm and the 1-norm on \mathbb{R}^n are comparable by showing that

$$\|x\|_\infty \leq \|x\|_1 \leq n\|x\|_\infty$$

for all $x \in \mathbb{R}^n$. Comparability of norms is an equivalence relation (Exercise 53), so it follows that the 1-norm and the 2-norm are comparable as well.

We now sketch the proof that *any vector norm $\|\cdot\|$ on \mathbb{R}^n is comparable to the 2-norm*. (A similar result holds for \mathbb{C}^n , but these results do *not* extend to infinite-dimensional spaces.) First we prove that there is a constant $C \in (0, \infty)$ such that for all $x = (x_1, \dots, x_n) \in \mathbb{R}^n$, $\|x\| \leq C\|x\|_2$. Let e_1, \dots, e_n be the standard basis vectors in \mathbb{R}^n , so that $x = x_1 e_1 + \dots + x_n e_n$. Let $C = \|e_1\| + \|e_2\| + \dots + \|e_n\|$, which is a positive finite constant. Now, use the norm axioms to compute

$$\begin{aligned} \|x\| &= \|x_1 e_1 + \dots + x_n e_n\| \leq \|x_1 e_1\| + \dots + \|x_n e_n\| \\ &= |x_1| \cdot \|e_1\| + \dots + |x_n| \cdot \|e_n\| \\ &\leq \|x\|_\infty \|e_1\| + \dots + \|x\|_\infty \|e_n\| = C\|x\|_\infty \leq C\|x\|_2. \end{aligned}$$

Obtaining the inequality $\|x\|_2 \leq D\|x\|$ requires a compactness argument. We will quote results from calculus concerning the metric space \mathbb{R}^n with the “Euclidean” metric $d(x, y) = \|x - y\|_2$. Let $S = \{x \in \mathbb{R}^n : \|x\|_2 = 1\}$ be the unit sphere in this metric space. The set S is closed and bounded in \mathbb{R}^n , hence is a compact subset of \mathbb{R}^n . Define $f : S \rightarrow \mathbb{R}$ by $f(x) = \|x\|$ for $x \in S$. We claim the function f is *continuous*. Given $\epsilon > 0$, let $\delta = \epsilon/C$. Then for $x, y \in S$ satisfying $d(x, y) = \|x - y\|_2 < \delta$, we have

$$|f(x) - f(y)| = |||x|| - ||y||| \leq \|x - y\| \leq C\|x - y\|_2 < \epsilon.$$

(The first inequality follows from Exercise 33.) From calculus, we know that a continuous real-valued function on a compact set attains its minimum value at some point in the set. So, there exists $x_0 \in S$ such that $f(x_0) \leq f(x)$ for all $x \in S$. Note that $x_0 \neq 0$ (since $\|x_0\|_2 = 1$), so $f(x_0) = \|x_0\| > 0$. Let $D = \|x_0\|^{-1}$, which is a positive finite constant. Given any $x \in \mathbb{R}^n$, we can now prove that $\|x\|_2 \leq D\|x\|$. This inequality certainly holds if $x = 0$. If $\|x\|_2 = 1$, then $x \in S$, so $f(x_0) \leq f(x)$, which means $D^{-1} \leq \|x\|$. Multiplying by D , we get $\|x\|_2 = 1 \leq D\|x\|$ as needed. Finally, consider an arbitrary nonzero $y \in \mathbb{R}^n$. Let $c = \|y\|_2$ and $x = c^{-1}y \in S$. Multiplying both sides of the known inequality $\|x\|_2 \leq D\|x\|$ by c , we get $\|cx\|_2 \leq D\|cx\|$, so $\|y\|_2 \leq D\|y\|$. This completes the proof that $\|\cdot\|$ and $\|\cdot\|_2$ are comparable. By Exercise 53, any two norms on \mathbb{R}^n (or \mathbb{C}^n) are comparable.

Finally, given comparable norms $\|\cdot\|$ and $\|\cdot\|'$, we show that $v_k \rightarrow v$ relative to the first norm iff $v_k \rightarrow v$ relative to the second norm. Choose constants C and D as in the definition of comparable norms. Suppose $v_k \rightarrow v$ relative to $\|\cdot\|$. Given $\epsilon > 0$, choose k_0 so that $k \geq k_0$ implies $\|v_k - v\| < \epsilon/D$. Then $k \geq k_0$ implies $\|v_k - v\|' \leq D\|v_k - v\| < \epsilon$. So $v_k \rightarrow v$ relative to $\|\cdot\|'$. Conversely, suppose $v_k \rightarrow v$ relative to $\|\cdot\|'$. Given $\epsilon > 0$, choose k_1 so that $k \geq k_1$ implies $\|v_k - v\|' < \epsilon/C$. Then $k \geq k_1$ implies $\|v_k - v\| \leq C\|v_k - v\|' < \epsilon$. Combining this result with the preceding theorem, we conclude that in \mathbb{R}^n (or \mathbb{C}^n), $v_k \rightarrow v$ relative to any given vector norm iff $v_k \rightarrow v$ relative to the 2-norm (or the 1-norm or the sup norm).

10.8 Matrix Norms

Recall that $M_n(\mathbb{C})$ is the set of all $n \times n$ matrices with entries in \mathbb{C} . Since $M_n(\mathbb{C})$ is a complex vector space, we can consider vector norms defined on $M_n(\mathbb{C})$. All such norms satisfy $0 \leq \|A\| < \infty$, $\|A\| = 0$ iff $A = 0$, $\|cA\| = |c|\cdot\|A\|$, and $\|A + B\| \leq \|A\| + \|B\|$ for all $A, B \in M_n(\mathbb{C})$ and all $c \in \mathbb{C}$. A *matrix norm* is a vector norm $\|\cdot\|$ on $M_n(\mathbb{C})$ satisfying the additional *submultiplicativity* axiom $\|AB\| \leq \|A\|\cdot\|B\|$ for all $A, B \in M_n(\mathbb{C})$.

For example, define $\|A\| = n \cdot \max_{i,j \in [n]} |A(i, j)|$ for $A \in M_n(\mathbb{C})$. The axioms for a vector norm can be checked as in §10.4. To check submultiplicativity, fix $A, B \in M_n(\mathbb{C})$ and $i, j \in [n]$. We know $(AB)(i, j) = \sum_{k=1}^n A(i, k)B(k, j)$, so

$$n|(AB)(i, j)| \leq \sum_{k=1}^n n|A(i, k)| \cdot |B(k, j)| \leq \sum_{k=1}^n n(n^{-1}\|A\|)(n^{-1}\|B\|) = \|A\|\cdot\|B\|.$$

This holds for all i, j , so taking the maximum gives $\|AB\| \leq \|A\|\cdot\|B\|$ as needed. (The inequality would not hold if we omitted the n from the definition of $\|A\|$.)

Given any vector norm $\|\cdot\|$ on \mathbb{C}^n , we now construct a matrix norm on $M_n(\mathbb{C})$ called the matrix norm *induced* by the given vector norm. In addition to the properties listed above, this matrix norm will satisfy $\|I_n\| = 1$ and $\|Av\| \leq \|A\|\cdot\|v\|$ for $A \in M_n(\mathbb{C})$ and $v \in \mathbb{C}^n$.

To define the matrix norm, let $U = \{v \in \mathbb{C}^n : \|v\| = 1\}$ be the set of unit vectors in \mathbb{C}^n relative to the given vector norm. For $A \in M_n(\mathbb{C})$, let

$$\|A\| = \sup\{\|Av\| : v \in U\}. \quad (10.8)$$

So the norm of A is the least upper bound of the set of lengths $\|Av\|$ as v ranges over all unit vectors in \mathbb{C}^n (we view elements of \mathbb{C}^n as column vectors here). For example, $\|0\| = \sup\{\|0v\| : v \in U\} = \sup\{0\} = 0$ and $\|I_n\| = \sup\{\|I_nv\| : v \in U\} = \sup\{1\} = 1$.

We now prove that our definition satisfies the requirements of a matrix norm. Fix A, B in $M_n(\mathbb{C})$ and $c \in \mathbb{C}$. For the first axiom, note that $0 \leq \|A\| \leq \infty$ since $\|A\|$ is the least upper bound of a set of nonnegative real numbers. Seeing that $\|A\|$ is *finite* is a bit tricky. We know $\|\cdot\|$ is comparable to $\|\cdot\|_\infty$, so there is a positive finite constant C with $\|v\|_\infty \leq C\|v\|$ for all $v \in \mathbb{C}^n$. Given any $v = (v_1, \dots, v_n) \in U$, we have $v = v_1e_1 + \dots + v_ne_n$, so $Av = A \sum_{i=1}^n v_i e_i = \sum_{i=1}^n v_i(Ae_i)$. Taking norms gives

$$\begin{aligned} \|Av\| &\leq \sum_{i=1}^n \|v_i(Ae_i)\| = \sum_{i=1}^n |v_i| \cdot \|Ae_i\| \\ &\leq \sum_{i=1}^n \|v\|_\infty \|Ae_i\| \leq \sum_{i=1}^n C\|v\| \cdot \|Ae_i\| = C \sum_{i=1}^n \|Ae_i\|. \end{aligned}$$

Thus, $C \sum_{i=1}^n \|Ae_i\|$ is a finite upper bound for the set $\{\|Av\| : v \in U\}$, so the least upper bound $\|A\|$ of this set is indeed finite. We already saw that $\|0\| = 0$. On the other hand, a nonzero matrix A must have some nonzero column $A^{[j]} = Ae_j$. Taking $v = \|e_j\|^{-1}e_j$, which is a unit vector, we have $\|Av\| = \|e_j\|^{-1}\|Ae_j\| \neq 0$, so $\|A\| > 0$.

For the second axiom, note that $\|(cA)v\| = |c| \cdot \|Av\|$ for all $v \in U$. So the least upper bound of the numbers $\|(cA)v\|$ over all $v \in U$ is $|c|$ times the least upper bound of the numbers $\|Av\|$, giving $\|cA\| = |c| \cdot \|A\|$. For the third axiom, note that $\|(A+B)v\| = \|Av + Bv\| \leq \|Av\| + \|Bv\| \leq \|A\| + \|B\|$ for all $v \in U$. Thus $\|A\| + \|B\|$ is an upper bound for the set $\{\|(A+B)v\| : v \in U\}$, and hence the least upper bound $\|A+B\|$ of this set is $\leq \|A\| + \|B\|$.

Next we check that $\|Av\| \leq \|A\| \cdot \|v\|$ for $A \in M_n(\mathbb{C})$ and $v \in \mathbb{C}^n$. Both sides are zero when $v = 0$. For nonzero v , write $v = cu$ where $c = \|v\|$ and $u \in U$. Then compute

$$\|Av\| = \|A(cu)\| = \|c(Au)\| = |c| \cdot \|Au\| \leq |c| \cdot \|A\| = \|A\| \cdot \|v\|.$$

For submultiplicativity, fix $A, B \in M_n(\mathbb{C})$. For any $v \in U$, the fact just proved gives

$$\|(AB)v\| = \|A(Bv)\| \leq \|A\| \cdot \|Bv\| \leq \|A\| \cdot \|B\| \cdot \|v\| = \|A\| \cdot \|B\|.$$

Thus $\|A\| \cdot \|B\|$ is an upper bound for the set $\{\|(AB)v\| : v \in U\}$, and hence the least upper bound $\|AB\|$ of this set is $\leq \|A\| \cdot \|B\|$.

Not all matrix norms arise from vector norms via (10.8). For example, the norm $\|A\| = n \cdot \max_{i,j} |A(i,j)|$ cannot arise from any vector norm since $\|I_n\| = n \neq 1$. Let us call a matrix norm $\|\cdot\|$ *induced* iff there exists a vector norm for which $\|\cdot\|$ is the associated matrix norm.

The proof of the following lemma uses the matrix norm induced by a given vector norm. Suppose $(v_k : k \geq 0)$ is a sequence in \mathbb{C}^n converging to $v \in \mathbb{C}^n$ relative to a vector norm $\|\cdot\|$; then for all $A \in M_n(\mathbb{C})$, $Av_k \rightarrow Av$. The result holds when $A = 0$, so assume $A \neq 0$. Fix $\epsilon > 0$, and choose $k \in \mathbb{N}$ such that $k \geq k_0$ implies $\|v_k - v\| < \epsilon/\|A\|$, where $\|A\|$ is the matrix norm induced by the given vector norm. Then $k \geq k_0$ implies $\|Av_k - Av\| = \|A(v_k - v)\| \leq \|A\| \cdot \|v_k - v\| < \epsilon$. For this proof to work, it is critical to know (as proved above) that $\|A\|$ is finite.

10.9 Formulas for Matrix Norms

In this section, we develop some explicit formulas for computing the matrix norms induced by the sup norm, the 1-norm, and the 2-norm on \mathbb{C}^n . We start by showing that the matrix norm $\|A\|_\infty$ induced by the vector norm $\|v\|_\infty$ satisfies

$$\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |A(i, j)|. \quad (10.9)$$

Let $U = \{v \in \mathbb{C}^n : \|v\|_\infty = 1\}$. For any $v \in U$, we compute

$$\begin{aligned} \|Av\|_\infty &= \max_{1 \leq i \leq n} |(Av)_i| = \max_{1 \leq i \leq n} \left| \sum_{j=1}^n A(i, j)v_j \right| \\ &\leq \max_{1 \leq i \leq n} \sum_{j=1}^n |A(i, j)| \cdot |v_j| \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |A(i, j)| \cdot \|v\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |A(i, j)|. \end{aligned}$$

So the expression on the right side of (10.9) is an upper bound for $\{\|Av\|_\infty : v \in U\}$, hence $\|A\|_\infty \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |A(i, j)|$. To establish the reverse inequality, fix an index $k \in [n]$ for which $\sum_{j=1}^n |A(k, j)|$ attains its maximum. For each $j \in [n]$, choose v_j to be a complex number of modulus 1 such that $A(k, j)v_j = |A(k, j)|$. (If $A(k, j) = re^{i\theta}$ in polar form, we can take $v_j = e^{-i\theta}$.) Now observe that $v = (v_1, \dots, v_n) \in \mathbb{C}^n$ has $\|v\|_\infty = 1$ since $|v_j| = 1$ for all j . So $v \in U$, and we conclude that

$$\begin{aligned} \|A\|_\infty &\geq \|Av\|_\infty = \max_{1 \leq i \leq n} |(Av)_i| \geq |(Av)_k| = \left| \sum_{j=1}^n A(k, j)v_j \right| = \sum_{j=1}^n |A(k, j)| \\ &= \max_{1 \leq i \leq n} \sum_{j=1}^n |A(i, j)|. \end{aligned}$$

Next we prove that the matrix norm $\|A\|_1$ induced by the vector norm $\|v\|_1$ satisfies

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |A(i, j)|. \quad (10.10)$$

We now take $U = \{v \in \mathbb{C}^n : \|v\|_1 = 1\}$. For any $v \in U$,

$$\begin{aligned} \|Av\|_1 &= \sum_{i=1}^n |(Av)_i| = \sum_{i=1}^n \left| \sum_{j=1}^n A(i, j)v_j \right| \\ &\leq \sum_{i=1}^n \sum_{j=1}^n |A(i, j)| \cdot |v_j| = \sum_{j=1}^n |v_j| \sum_{i=1}^n |A(i, j)| \\ &\leq \sum_{j=1}^n |v_j| \left[\max_{1 \leq j \leq n} \sum_{i=1}^n |A(i, j)| \right] = \|v\|_1 \max_{1 \leq j \leq n} \sum_{i=1}^n |A(i, j)| = \max_{1 \leq j \leq n} \sum_{i=1}^n |A(i, j)|. \end{aligned}$$

We conclude as before that $\|A\|_1 \leq \max_{1 \leq j \leq n} \sum_{i=1}^n |A(i, j)|$. For the reverse inequality, fix

an index $k \in [n]$ for which $\sum_{i=1}^n |A(i, k)|$ is maximized. Since $\|e_k\|_1 = 1$, $e_k \in U$, and hence

$$\|A\|_1 \geq \|Ae_k\|_1 = \|A^{[k]}\|_1 = \sum_{i=1}^n |A(i, k)| = \max_{1 \leq j \leq n} \sum_{i=1}^n |A(i, j)|.$$

Let $\|A\|_2$ be the matrix norm induced by the vector norm $\|v\|_2$. For any $B \in M_n(\mathbb{C})$, define the *spectral radius* $\rho(B)$ to be the maximum magnitude of all the complex eigenvalues of B , i.e.,

$$\rho(B) = \max\{|c| : c \in \mathbb{C} \text{ is an eigenvalue of } B\}.$$

It can be shown (Exercise 60) that $\|A\|_2 = \sqrt{\rho(A^*A)}$ for $A \in M_n(\mathbb{C})$. One can also show that $\|A\|_2$ is the largest singular value of A (see Exercise 55 in Chapter 9).

10.10 Matrix Inversion via Geometric Series

In calculus, one learns the *geometric series* formula

$$\frac{1}{1-r} = 1 + r + r^2 + r^3 + \cdots + r^k + \cdots = \sum_{k=0}^{\infty} r^k,$$

which is valid for all complex numbers r such that $|r| < 1$. In this section, we prove an analogous formula in which r is replaced by an $n \times n$ complex matrix C . We will assume that $\|C\| < 1$ for some induced matrix norm on \mathbb{C}^n ; this is the analogue of the hypothesis $|r| < 1$.

We intend to show that $I - C$ is an invertible matrix, and

$$(I - C)^{-1} = I + C + C^2 + \cdots + C^k + \cdots = \sum_{k=0}^{\infty} C^k, \quad (10.11)$$

where the infinite series of matrices denotes the limit of the sequence of partial sums $(I + C + C^2 + \cdots + C^k : k \geq 0)$. To get a contradiction, assume that $I - C$ is not invertible. Then there exists a nonzero vector x with $(I - C)x = 0$, so $x = Cx$. Taking norms gives $\|x\| = \|Cx\| \leq \|C\| \cdot \|x\|$. Dividing by the positive real number $\|x\|$ gives $1 \leq \|C\|$, contradicting our hypothesis on C . So $B = (I - C)^{-1}$ does exist. B must be nonzero, so $\|B\| > 0$.

Now fix $k \geq 0$, and compute

$$(I - C)(I + C + C^2 + \cdots + C^k) = (I + C + C^2 + \cdots + C^k) - (C + C^2 + \cdots + C^k + C^{k+1}) = I - C^{k+1}.$$

Multiplying both sides by B gives $I + C + C^2 + \cdots + C^k = B - BC^{k+1}$. Given $\epsilon > 0$, we can choose k_0 so that $k \geq k_0$ implies $\|C\|^{k+1} < \epsilon/\|B\|$; this is possible since $\|C\| < 1$. For any $k \geq k_0$, we then have

$$\|(I + C + C^2 + \cdots + C^k) - B\| = \|-BC^{k+1}\| \leq \|B\| \cdot \|C\|^{k+1} < \epsilon.$$

This proves that the sequence of partial sums converges to $B = (I - C)^{-1}$, as needed.

10.11 Affine Iteration and Richardson's Algorithm

In this section, we study the convergence of the iterative algorithm obtained by repeatedly applying an affine map to a given initial vector. The input to the algorithm consists of a fixed matrix $C \in M_n(\mathbb{C})$, a vector $d \in \mathbb{C}^n$, and an initial vector $x_0 \in \mathbb{C}^n$. We construct a sequence $(x_k : k \geq 0)$ in \mathbb{C}^n by the recursive formula

$$x_{k+1} = Cx_k + d \quad (k \geq 0), \quad (10.12)$$

and we seek conditions under which this sequence of vectors will converge to some limit $x \in \mathbb{C}^n$.

If the limit x exists at all, we can find it by letting k go to infinity on both sides of (10.12). Assuming that $x_k \rightarrow x$, the left side x_{k+1} will also converge to x . On the other side, Cx_k will converge to Cx , so $Cx_k + d$ will converge to $Cx + d$. Then $x = Cx + d$, so that the limit x must solve the linear system $(I - C)x = d$. If $I - C$ is invertible, then the only possible limit of the iterative algorithm is $x = (I - C)^{-1}d$.

To prove convergence, we will assume that $\|C\| < 1$ for some induced matrix norm on \mathbb{C}^n . In this case, we saw in the last section that $I - C$ is invertible. We will prove that for any choice of initial vector $x_0 \in \mathbb{C}^n$, the sequence defined by (10.12) converges to $x = (I - C)^{-1}d$, and we will derive a bound on the error $\|x_k - x\|$. We know $x_{k+1} = Cx_k + d$ for $k \geq 0$, and $x = Cx + d$. Subtracting these equations gives $x_{k+1} - x = (Cx_k + d) - (Cx + d) = Cx_k - Cx = C(x_k - x)$. Taking norms, we get

$$\|x_{k+1} - x\| = \|C(x_k - x)\| \leq \|C\| \cdot \|x_k - x\|$$

for all $k \geq 0$. If $k \geq 1$, we can write $\|x_k - x\| \leq \|C\| \cdot \|x_{k-1} - x\|$ on the right side here, giving $\|x_{k+1} - x\| \leq \|C\|^2 \|x_{k-1} - x\|$. Continuing in this way, we eventually obtain the bound

$$\|x_{k+1} - x\| \leq \|C\|^{k+1} \|x_0 - x\|,$$

valid for all $k \geq 0$. The right side goes to zero at an exponential rate as k goes to infinity, since $\|C\| < 1$. Therefore, $\lim_{k \rightarrow \infty} \|x_k - x\| = 0$, proving that $x_k \rightarrow x$ as needed.

This iterative method finds the unique x solving $(I - C)x = d$. So, if we want to use this method to solve a linear system $Ax = b$ (for given $A \in M_n(\mathbb{C})$ and $b \in \mathbb{C}^n$), we can choose $C = I - A$ and $d = b$. Provided that $r = \|I - A\| < 1$ for some induced matrix norm, the method will converge to the solution x , with the error in the k 'th term bounded by r^k times the initial error $\|x_0 - x\|$. Since $C = I - A$, the iteration formula (10.12) becomes

$$x_{k+1} = x_k + (b - Ax_k) \quad (k \geq 0),$$

so we have recovered Richardson's algorithm for solving a linear system.

We can use the formulas for matrix norms in §10.9 to find explicit sufficient conditions on A that will guarantee that Richardson's algorithm will converge. First, if $\max_{1 \leq i \leq n} \sum_{j=1}^n |(I - A)(i, j)| < 1$, then (10.9) shows that $\|I - A\|_\infty < 1$, so the algorithm converges. Second, if $\max_{1 \leq j \leq n} \sum_{i=1}^n |(I - A)(i, j)| < 1$, then (10.10) shows that $\|I - A\|_1 < 1$, so the algorithm converges. For example, the matrix $A = \begin{bmatrix} 1.1 & 0.3 \\ -0.4 & 0.8 \end{bmatrix}$ considered in §10.1 satisfies both conditions just mentioned, so Richardson's algorithm for solving $Ax = b$ will converge for any choice of b and x_0 . In fact, since $\|I - A\|_1 = 1/2$, we know that $\|x_k - x\| \leq 2^{-k} \|x_0 - x\|$ for all $k \geq 1$, so that the bound for the maximum error is cut in half for each additional iteration.

10.12 Splitting Matrices and Jacobi's Algorithm

We now apply the analysis in the preceding section to iterative algorithms that solve $Ax = b$ by the following procedure. The algorithm picks an invertible *splitting matrix* Q (depending on the given $A \in M_n(\mathbb{C})$) and computes $(x_k : k \geq 0)$ via the iteration formula

$$Qx_{k+1} = b - (A - Q)x_k \quad (k \geq 0). \quad (10.13)$$

For example, Richardson's algorithm takes $Q = I$; Jacobi's algorithm takes Q to be the diagonal part of A ; and the Gauss–Seidel algorithm takes Q to be the lower-triangular part of A . Recall that an algorithm of this kind will only be practical if $Qx_{k+1} = y$ can be solved quickly for x_{k+1} given y .

By multiplying both sides of (10.13) on the left by Q^{-1} , we see that (10.13) is algebraically (though not computationally) equivalent to

$$x_{k+1} = Q^{-1}b - (Q^{-1}A - Q^{-1}Q)x_k = (I - Q^{-1}A)x_k + Q^{-1}b.$$

Therefore, we can invoke the results of §10.11 taking $C = I - Q^{-1}A$ and $d = Q^{-1}b$. Note that $(I - C)x = d$ iff $Q^{-1}Ax = Q^{-1}b$ iff $Ax = b$ for this choice of C and d . We conclude that if $\|I - Q^{-1}A\| < 1$ for some induced matrix norm, then the sequence defined by (10.13) will converge to the true solution for any starting vector x_0 . Moreover, we have the error bound

$$\|x_k - x\| \leq \|I - Q^{-1}A\|^k \|x_0 - x\| \quad (10.14)$$

for all $k \geq 0$.

Using this analysis, we can prove a sufficient condition on A that guarantees the convergence of the Jacobi method for solving $Ax = b$. Call $A \in M_n(\mathbb{C})$ *diagonally dominant* iff

$$|A(i, i)| > \sum_{\substack{j=1 \\ j \neq i}}^n |A(i, j)| \text{ for all } i \in [n].$$

(Such a matrix must have all diagonal entries nonzero.) We now show that *the Jacobi method always converges given a diagonally dominant input matrix A*. In this case, Q^{-1} is a diagonal matrix with $Q^{-1}(i, i) = A(i, i)^{-1}$. Therefore, the i, j -entry of $I - Q^{-1}A$ is 0 if $i = j$, and $-A(i, i)^{-1}A(i, j)$ otherwise. So for each $i \in [n]$,

$$\sum_{j=1}^n |(I - Q^{-1}A)(i, j)| = \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|A(i, j)|}{|A(i, i)|} < 1.$$

Taking the maximum over $i \in [n]$, it follows from (10.9) that $\|I - Q^{-1}A\|_\infty < 1$, proving convergence.

10.13 Induced Matrix Norms and the Spectral Radius

In the previous sections, we proved that various iterative algorithms converge provided that a certain matrix $B \in M_n(\mathbb{C})$ satisfied $\|B\| < 1$ for some induced matrix norm. It is possible that $\|B\|_1 \geq 1$, $\|B\|_\infty \geq 1$, and yet $\|B\| < 1$ for some other matrix norm induced by some

vector norm other than $\|\cdot\|_1$ or $\|\cdot\|_\infty$. To obtain a sufficient condition for convergence that is as powerful as possible, we would really like to know the quantity

$$f(B) = \inf\{\|B\| : \|\cdot\| \text{ is an induced matrix norm on } \mathbb{C}^n\}, \quad (10.15)$$

which is the greatest lower bound in \mathbb{R} of the numbers $\|B\|$ as the norm varies over all possible induced matrix norms. Recall from §10.9 that the spectral radius of B is

$$\rho(B) = \max\{|c| : c \in \mathbb{C} \text{ is an eigenvalue of } B\}.$$

We will prove that $f(B) = \rho(B)$ for all $B \in M_n(\mathbb{C})$.

To begin, let $c \in \mathbb{C}$ be an eigenvalue of B with $|c| = \rho(B)$, and let $v \neq 0$ be an associated eigenvector. Fix an arbitrary vector norm $\|\cdot\|$ on \mathbb{C}^n , and note that $u = \|v\|^{-1}v$ satisfies $\|u\| = 1$ and $Bu = cu$. Therefore, using the matrix norm induced by this vector norm, we compute

$$\|B\| \geq \|Bu\| = \|cu\| = |c| \cdot \|u\| = |c| = \rho(B).$$

So $\rho(B)$ is a lower bound for the set of numbers on the right side of (10.15). As $f(B)$ is the greatest lower bound of this set, $\rho(B) \leq f(B)$ follows.

Showing the opposite inequality $f(B) \leq \rho(B)$ requires more work. It is enough to prove that $f(B) \leq \rho(B) + \epsilon$ for each fixed $\epsilon > 0$. Given $\epsilon > 0$, we first show that there is an invertible matrix $S \in M_n(\mathbb{C})$ such that $S^{-1}BS$ is upper-triangular and $\|S^{-1}BS\|_\infty \leq \rho(B) + \epsilon$. We give an argument using Jordan canonical forms (Chapter 8); a different argument based on unitary triangularization of complex matrices is sketched in Exercise 61. Define $T : \mathbb{C}^n \rightarrow \mathbb{C}^n$ by $T(v) = Bv$ for $v \in \mathbb{C}^n$. We know there is an ordered basis $X = (x_1, x_2, \dots, x_n)$ for \mathbb{C}^n such that $[T]_X$ is a Jordan canonical form. This means that for certain scalars $c_1, \dots, c_n \in \mathbb{C}$ and $d_1, \dots, d_n \in \{0, 1\}$, $T(x_1) = c_1x_1$ and $T(x_i) = c_i x_i + d_i x_{i-1}$ for $2 \leq i \leq n$. Now, replace the ordered basis X by the ordered basis $Y = (y_1, y_2, \dots, y_n)$, where $y_i = \epsilon^i x_i$ for all $i \in [n]$. We compute $T(y_1) = c_1 y_1$ and

$$T(y_i) = T(\epsilon^i x_i) = \epsilon^i T(x_i) = \epsilon^i(c_i x_i + d_i x_{i-1}) = c_i y_i + (d_i \epsilon) y_{i-1}$$

for $2 \leq i \leq n$. Letting $S = {}_E[\text{id}]_Y$, where E is the standard ordered basis of \mathbb{C}^n , we see that $U = S^{-1}BS = [T]_Y$ is an upper-triangular matrix with main diagonal entries c_1, \dots, c_n , entries equal to zero or ϵ on the next higher diagonal, and zeroes elsewhere. Since U is triangular, the eigenvalues of U are c_1, \dots, c_n . Since B is similar to U , these are also the eigenvalues of B . For each $i \in [n]$, $\sum_{j=1}^n |U(i, j)|$ is either $|c_i|$ or $|c_i| + \epsilon$. By (10.9), we see that

$$\|U\|_\infty \leq \max_{1 \leq i \leq n} (|c_i| + \epsilon) = \rho(B) + \epsilon.$$

To continue, we need a clever choice of a vector norm and its induced matrix norm. In Exercise 34, we ask the reader to check that for any fixed invertible $S \in M_n(\mathbb{C})$, the formula $\|v\|_S = \|S^{-1}v\|_\infty$ for $v \in \mathbb{C}^n$ defines a vector norm on \mathbb{C}^n . Moreover, the induced matrix norm satisfies $\|A\|_S = \|S^{-1}AS\|_\infty$ for all $A \in M_n(\mathbb{C})$. Let us choose S to be the matrix found in the previous paragraph. Then

$$f(B) \leq \|B\|_S = \|S^{-1}BS\|_\infty = \|U\|_\infty \leq \rho(B) + \epsilon,$$

as needed.

We have now proved that $f(B) = \rho(B)$ for all $B \in M_n(\mathbb{C})$. Using (10.15), we see that $\rho(B) < 1$ iff there exists an induced matrix norm such that $\|B\| < 1$. Accordingly, we can restate the convergence result in §10.12 as follows. *An iterative algorithm based on (10.13) will converge if the spectral radius $\rho(I - Q^{-1}A)$ is less than 1.*

10.14 Analysis of the Gauss–Seidel Algorithm

This section uses the convergence criterion $\rho(I - Q^{-1}A) < 1$ to prove that *the Gauss–Seidel algorithm will converge if the input matrix $A \in M_n(\mathbb{C})$ is diagonally dominant*. Recall that diagonal dominance means $|A(i, i)| > \sum_{j \neq i} |A(i, j)|$ for all $i \in [n]$. Now, let $c \in \mathbb{C}$ be any eigenvalue of $I - Q^{-1}A$ with associated eigenvector $x = (x_1, \dots, x_n) \neq 0$. It will suffice to prove $|c| < 1$.

We know $(I - Q^{-1}A)x = cx$. Left-multiplying both sides by Q gives $(Q - A)x = c(Qx)$. Recall that, in the Gauss–Seidel algorithm, Q contains the entries of A on or below the main diagonal, and $Q - A$ involves the entries of A above the main diagonal. So, taking the i 'th component of $(Q - A)x = c(Qx)$ gives

$$-\sum_{j:j>i} A(i, j)x_j = c \sum_{j:j \leq i} A(i, j)x_j$$

for all $i \in [n]$. Isolating the term $cA(i, i)x_i$ and taking magnitudes gives

$$|c| |A(i, i)| |x_i| = \left| -\sum_{j:j>i} A(i, j)x_j - c \sum_{j:j < i} A(i, j)x_j \right| \leq \sum_{j:j>i} |A(i, j)| |x_j| + |c| \sum_{j:j < i} |A(i, j)| |x_j|.$$

Fix an index i such that $|x_i| = \max_{1 \leq k \leq n} |x_k| > 0$. Then $|x_j|/|x_i| \leq 1$ for all $j \neq i$, so dividing the preceding inequality by $|x_i|$ gives

$$|c| |A(i, i)| \leq \sum_{j:j>i} |A(i, j)| + |c| \sum_{j:j < i} |A(i, j)|. \quad (10.16)$$

Now, the diagonal dominance of A implies that $|A(i, i)| - \sum_{j:j < i} |A(i, j)| > \sum_{j:j > i} |A(i, j)| \geq 0$. Therefore, we can solve (10.16) for $|c|$ to conclude that

$$|c| \leq \frac{\sum_{j:j>i} |A(i, j)|}{|A(i, i)| - \sum_{j:j < i} |A(i, j)|} < 1,$$

as needed.

10.15 Power Method for Finding Eigenvalues

To compute the spectral radius of $A \in M_n(\mathbb{C})$ from the definition, we need to know the largest complex eigenvalue of A . We now discuss an iterative algorithm called the *power method* whose goal is to compute this largest eigenvalue and an associated eigenvector. The algorithm takes as input the matrix A and an arbitrary nonzero initial vector $x_0 \in \mathbb{C}^n$. The algorithm iteratively computes $y_{k+1} = Ax_k$ and $x_{k+1} = y_{k+1}/\|y_{k+1}\|_\infty$ for all $k \geq 0$. (If we get $y_{k+1} = 0$, the algorithm fails, and we try again with a different x_0 .) At stage k , x_k is the algorithm's approximation for the required eigenvector. The associated eigenvalue is estimated by choosing an index $i \in [n]$ and returning $c_k = (Ax_k)(i)/x_k(i)$. Any index i can be used here, as long as $|x_k(i)|$ is not too close to zero.

For example, let $A = \begin{bmatrix} 1 & 4 & -1 \\ 0 & 3 & 2 \\ 1 & -1 & -3 \end{bmatrix}$ and $x_0 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$. The vectors x_1, x_2, \dots, x_{12} are:

$$\begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0.2 \\ 0.4 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ -0.286 \end{bmatrix}, \begin{bmatrix} 1 \\ 0.459 \\ 0.162 \end{bmatrix}, \begin{bmatrix} 1 \\ 0.636 \\ 0.020 \end{bmatrix}, \\ \begin{bmatrix} 1 \\ 0.553 \\ 0.086 \end{bmatrix}, \begin{bmatrix} 1 \\ 0.586 \\ 0.060 \end{bmatrix}, \begin{bmatrix} 1 \\ 0.572 \\ 0.071 \end{bmatrix}, \begin{bmatrix} 1 \\ 0.578 \\ 0.067 \end{bmatrix}, \begin{bmatrix} 1 \\ 0.575 \\ 0.068 \end{bmatrix}, \begin{bmatrix} 1 \\ 0.576 \\ 0.068 \end{bmatrix}.$$

At this stage, our estimate for the eigenvector is $x_{12} = (1, 0.576, 0.068)$, and $Ax_{12} = (3.236, 1.864, 0.22)$. Dividing each entry of Ax_{12} by the corresponding entry of x_{12} gives us three estimates for the largest eigenvalue of A : 3.236 , $1.864/0.576 \approx 3.236$, and $0.22/0.068 \approx 3.235$. In fact, by calculating the characteristic polynomial of A , one sees that the exact eigenvalues of A are -1 , $1 - \sqrt{5}$, and $1 + \sqrt{5}$. Thus, $\rho(A) = 1 + \sqrt{5} \approx 3.23607$, so that we obtain quite good estimates for $\rho(A)$ after twelve iterations of the algorithm.

We now discuss a sufficient condition on A guaranteeing that the approximations produced by the power method will converge to $\rho(A)$. The condition we impose on A is that A is a *diagonalizable* matrix having a *unique* eigenvalue $c \in \mathbb{C}$ such that $|c| = \rho(A)$. (Some of the exercises explore what happens if these conditions are not met.) Given these assumptions, we can choose an ordered basis $Z = (z_1, \dots, z_n)$ of \mathbb{C}^n consisting of n linearly independent eigenvectors of A . We have $Az_i = c_i z_i$ for the complex eigenvalues c_1, \dots, c_n of A . We can reorder the basis Z so that

$$\rho(A) = |c_1| > |c_2| \geq |c_3| \geq \dots \geq |c_n|.$$

Given any $x_0 \in \mathbb{C}^n$, we obtain x_k by a sequence of k steps, each of which involves multiplying by A and then rescaling to get a unit vector. Letting $b_k \in \mathbb{R}^+$ be the product of the rescaling factors used to reach x_k , we see that $x_k = b_k A^k x_0$ for all $k \geq 1$.

Write $x_0 = d_1 z_1 + d_2 z_2 + \dots + d_n z_n$ for unique scalars $d_i \in \mathbb{C}$. We assume $d_1 \neq 0$, which will almost surely happen if x_0 is chosen at random in \mathbb{C}^n . By replacing each z_i in the basis Z by $d_i z_i$, we can assume without loss of generality that every d_i is 1. We have $A^k z_i = c_i^k z_i$ for all $i \in [n]$ and all $k \geq 1$, so $x_k = b_k(c_1^k z_1 + c_2^k z_2 + \dots + c_n^k z_n)$. Letting $b'_k = b_k c_1^k$ and $r_i = c_i/c_1$ for $2 \leq i \leq n$, we can rewrite this as

$$x_k = b'_k(z_1 + r_2^k z_2 + \dots + r_n^k z_n).$$

We have $|r_i| = |c_i|/|c_1| < 1$ for $2 \leq i \leq n$, so for each such i , the sequence $r_i^k z_i$ goes to zero as k goes to infinity. Letting $x'_k = x_k/b'_k$, it follows that $x'_k = z_1 + r_2^k z_2 + \dots + r_n^k z_n \rightarrow z_1$ as $k \rightarrow \infty$.

Next, let $i \in [n]$ be any index such that $z_1(i) \neq 0$. Since $x'_k \rightarrow z_1$ relative to $\|\cdot\|_\infty$, the inequality $0 \leq |x'_k(i) - z_1(i)| \leq \|x'_k - z_1\|_\infty$ shows that $x'_k(i) \rightarrow z_1(i)$ in \mathbb{C} as k goes to infinity. Similarly, since $Ax'_k \rightarrow Az_1 = c_1 z_1$, we see that $(Ax'_k)(i) \rightarrow c_1 z_1(i)$ as k goes to infinity. Finally, since $x_k = b'_k x'_k$,

$$\lim_{k \rightarrow \infty} \frac{(Ax_k)(i)}{x_k(i)} = \lim_{k \rightarrow \infty} \frac{(Ax'_k)(i)}{x'_k(i)} = \frac{c_1 z_1(i)}{z_1(i)} = c_1,$$

which says that the approximations produced by the power method do converge to the largest eigenvalue c_1 of A .

Although x'_k converges to the eigenvector z_1 associated with the eigenvalue c_1 , we cannot

conclude that the sequence x_k itself converges to any fixed eigenvector of A . The trouble is that the complex scaling constant $b'_k = b_k c_1^k$ may cause the sequence x_k to “jump around” between several eigenvectors that differ by a complex scaling factor of modulus 1. For example, if $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ and x_0 is the eigenvector $\begin{bmatrix} i \\ 1 \end{bmatrix}$ associated with the eigenvalue i , the power method produces the periodic sequence

$$\begin{bmatrix} i \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ i \end{bmatrix}, \begin{bmatrix} -i \\ -1 \end{bmatrix}, \begin{bmatrix} 1 \\ -i \end{bmatrix}, \begin{bmatrix} i \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ i \end{bmatrix}, \begin{bmatrix} -i \\ -1 \end{bmatrix}, \begin{bmatrix} 1 \\ -i \end{bmatrix}, \dots$$

This sequence does not converge, but each vector in the sequence is an eigenvector associated with the eigenvalue i . On the other hand, one can adapt the analysis given above to show that $x_k - b'_k z_1 \rightarrow 0$ as $k \rightarrow \infty$ (Exercise 63). Informally, this means that for large enough k , x_k is guaranteed to come arbitrarily close to some eigenvector associated with the eigenvalue c_1 , but this eigenvector may depend on k . Alternatively, since b_k is known and c_1 is being estimated by the algorithm, we can use an estimate for c_1 to approximate $z_1 = x_k / (b_k c_1^k)$.

10.16 Shifted and Inverse Power Method

The power method provides an iterative algorithm for computing the largest eigenvalue of $A \in M_n(\mathbb{C})$. What if we want to find one of the other eigenvalues of A ? This section and the next one discuss several approaches to this question.

Observe first that the largest eigenvalue of A is the eigenvalue farthest from the origin in the complex plane. Suppose we replace the matrix A by a *shifted* matrix $A + bI$, where $b \in \mathbb{C}$ is a fixed constant. Let $c \in \mathbb{C}$ be an eigenvalue of A with associated eigenvector $x \in \mathbb{C}^n$. Since $Ax = cx$, we have $(A + bI)x = Ax + bx = (c + b)x$, so that $c + b \in \mathbb{C}$ is an eigenvalue of $A + bI$ with associated eigenvector x . Conversely, if c' is an eigenvalue of $A + bI$, a similar computation shows that $c' - b$ is an eigenvalue of A . Geometrically, the eigenvalues of $A + bI$ are obtained by shifting (or translating) all eigenvalues of A by the fixed scalar b . For any $c \in \mathbb{C}$, the distance from c to 0 in \mathbb{C} is the same as the distance from $c + b$ to b in \mathbb{C} . So the eigenvalue of A farthest from the origin gets shifted to the eigenvalue of $A + bI$ farthest from the point b .

Suppose $C \in M_n(\mathbb{C})$ is a given matrix, $b \in \mathbb{C}$ is a given scalar, and we want to compute the complex eigenvalue of C farthest from b . Taking $A = C - bI$ in the previous paragraph, we first compute the eigenvalue c of A farthest from the origin using the power method. Then the required eigenvalue for C is $c + b$. This algorithm is called the *shifted power method*.

Now suppose $A \in M_n(\mathbb{C})$ and we want to find the eigenvalue of A *closest* to the origin. If A is not invertible, this eigenvalue is zero. If A^{-1} exists and c is any eigenvalue of A with associated eigenvector x , then $Ax = cx$ implies $A^{-1}Ax = A^{-1}cx$, hence $x = cA^{-1}x$. As c must be nonzero, we see $A^{-1}x = c^{-1}x$, so that c^{-1} is an eigenvalue of A^{-1} with associated eigenvector x . The argument is reversible, so the eigenvalues of A^{-1} are the multiplicative inverses in \mathbb{C} of the eigenvalues of A . The inverse of a complex number written in polar form is $(re^{i\theta})^{-1} = r^{-1}e^{-i\theta}$. It follows that for nonzero $c, d \in \mathbb{C}$, c is closer to the origin than d iff c^{-1} is farther from the origin than d^{-1} . Using the power method to find the eigenvalue c of A^{-1} that is farthest from the origin, it follows that c^{-1} is the eigenvalue of A closest to the origin. This algorithm is called the *inverse power method*.

Finally, we can use the *shifted inverse power method* to find the eigenvalue of a given $C \in M_n(\mathbb{C})$ that is closest to a given $b \in M_n(\mathbb{C})$. We apply the original power method to the matrix $(C - bI)^{-1}$ to obtain the largest eigenvalue c , and then return $c^{-1} + b$.

When implementing the inverse power method (or its shifted version), one needs to compute $y_{k+1} = A^{-1}x_k$ given x_k . Equivalently, one must solve $Ay_{k+1} = x_k$ for the unknown vector y_{k+1} . In practice, it is usually more efficient to obtain y_{k+1} by using Gaussian elimination algorithms to solve $Ay_{k+1} = x_k$ rather than computing A^{-1} explicitly. This is especially true when A has special structure such as sparsity.

10.17 Deflation

Suppose $A \in M_n(\mathbb{C})$ is a given matrix and we have found, by whatever method, one particular eigenvalue c of A and an associated eigenvector x . This section describes a general algorithm called *deflation* that produces a new matrix $B \in M_{n-1}(\mathbb{C})$ whose eigenvalues (counted by their algebraic multiplicities as roots of the characteristic polynomial) are precisely the eigenvalues of A with one copy of c removed. In particular, by repeatedly applying the power method (or its variants) followed by deflation to a given matrix A , we can eventually compute all the eigenvalues of A .

Recall Schur's theorem from §7.7: for any matrix $A \in M_n(\mathbb{C})$, there exists a unitary matrix $U \in M_n(\mathbb{C})$ such that $U^{-1}AU = U^*AU$ is upper-triangular. We obtain the deflation algorithm by examining the computations implicit in the proof of this result. Assume $c_1 \in \mathbb{C}$ is the known eigenvalue of A with associated known eigenvector $x_1 \in \mathbb{C}^n$. Dividing x_1 by $\|x_1\|_2$, we can assume that $\|x_1\|_2 = 1$. Let $T : \mathbb{C}^n \rightarrow \mathbb{C}^n$ be defined by $T(v) = Av$ for $v \in \mathbb{C}^n$. The key step in the inductive proof of Schur's theorem was to extend the list (x_1) to an orthonormal basis $X = (x_1, x_2, \dots, x_n)$ of \mathbb{C}^n . We can compute such an orthonormal basis X explicitly using the Gram–Schmidt orthonormalization algorithm (see §9.2 and Exercise 12 in Chapter 9).

Arrange the column vectors x_1, x_2, \dots, x_n as the columns of a matrix U , so $U = {}_E[\text{id}]_X$ where $E = (e_1, \dots, e_n)$ is the standard ordered basis of \mathbb{C}^n . Since the columns of U are orthonormal, U is a unitary matrix (see §7.5), so $U^{-1} = U^*$. We can compute $C = [T]_X = U^*AU$ explicitly. The first column of this matrix is $(c_1, 0, \dots, 0)^T$. Let $B \in M_{n-1}(\mathbb{C})$ be obtained from C by deleting the first row and first column. Computing the characteristic polynomial of C by expanding the determinant down the first column, we see that $\chi_A(t) = \chi_C(t) = (t - c_1)\chi_B(t)$. Therefore, the matrix B has the required eigenvalues, and the deflation algorithm returns B as its output.

10.18 Summary

- Vector Norms.* Given a real or complex vector space V , a vector norm $\|\cdot\| : V \rightarrow \mathbb{R}$ satisfies $0 \leq \|x\| < \infty$, $\|x\| = 0$ iff $x = 0$, $\|cx\| = |c| \cdot \|x\|$, and $\|x + y\| \leq \|x\| + \|y\|$ for all $x, y \in V$ and all scalars c . The pair $(V, \|\cdot\|)$ is called a normed vector space. Examples of vector norms on \mathbb{R}^n and \mathbb{C}^n include the sup norm $\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$, the 1-norm $\|x\|_1 = \sum_{i=1}^n |x_i|$, and the 2-norm $\|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2}$.
- Metric Spaces.* A metric space is a set X and a metric $d : X \times X \rightarrow \mathbb{R}$ satisfying $0 \leq d(x, y) < \infty$, $d(x, y) = 0$ iff $x = y$, $d(x, y) = d(y, x)$, and $d(x, z) \leq d(x, y) + d(y, z)$ for all $x, y, z \in X$. A normed vector space V has

the metric $d(x, y) = \|x - y\|$, which also satisfies $d(x + z, y + z) = d(x, y)$ and $d(cx, cy) = |c|d(x, y)$ for all $x, y, z \in V$ and all scalars c .

3. *Convergence Properties.* In a metric space (X, d) , a sequence $(x_k : k \geq 0)$ converges to $x \in X$ iff for all $\epsilon > 0$, there exists $k_0 \in \mathbb{N}$ such that for all $k \geq k_0$, $d(x_k, x) < \epsilon$. The limit of a convergent sequence is unique when it exists. In \mathbb{C}^n , suppose $v_k \rightarrow v$ and $w_k \rightarrow w$. Then $v_k + w_k \rightarrow v + w$, $cv_k \rightarrow cv$ for any scalar c , and $Av_k \rightarrow Av$ for any matrix A .
4. *Comparable Norms.* Vector norms $\|\cdot\|$ and $\|\cdot\|'$ on a vector space V are comparable iff there exist real C, D with $0 < C, D < \infty$ such that for all $x \in V$, $\|x\| \leq C\|x\|'$ and $\|x\|' \leq D\|x\|$. Any two vector norms on \mathbb{R}^n or \mathbb{C}^n are comparable. In particular, $\|x\|_\infty \leq \|x\|_2 \leq \sqrt{n}\|x\|_\infty$ and $\|x\|_\infty \leq \|x\|_1 \leq n\|x\|_\infty$.
5. *Matrix Norms.* A matrix norm $\|\cdot\| : M_n(\mathbb{C}) \rightarrow \mathbb{R}$ satisfies $0 \leq \|A\| < \infty$, $\|A\| = 0$ iff $A = 0$, $\|cA\| = |c|\cdot\|A\|$, $\|A+B\| \leq \|A\| + \|B\|$, and $\|AB\| \leq \|A\|\cdot\|B\|$ for all $A, B \in M_n(\mathbb{C})$ and all $c \in \mathbb{C}$. Any vector norm $\|\cdot\|$ on \mathbb{C}^n has an associated matrix norm $\|A\| = \sup\{\|Av\| : v \in \mathbb{C}^n, \|v\| = 1\}$, which satisfies $\|Ax\| \leq \|A\|\cdot\|x\|$ for all $x \in \mathbb{C}^n$. A matrix norm arising in this way from some vector norm is called an induced matrix norm.
6. *Formulas for Matrix Norms.* Fix $A \in M_n(\mathbb{C})$. The matrix norm induced by the sup norm is $\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |A(i, j)|$. The matrix norm induced by the 1-norm is $\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |A(i, j)|$. The spectral radius $\rho(A)$ is the maximum of the numbers $|c|$ as c ranges over all complex eigenvalues of A . The spectral radius of A also equals $\inf\{\|A\|\}$, where $\|\cdot\|$ varies over all induced matrix norms on \mathbb{C}^n . The matrix norm induced by the 2-norm is $\|A\|_2 = \sqrt{\rho(A^*A)}$.
7. *Geometric Series for Matrix Inverses.* If $C \in M_n(\mathbb{C})$ satisfies $\|C\| < 1$ for some induced matrix norm (i.e., $\rho(C) < 1$), then $I - C$ is invertible and $(I - C)^{-1} = \sum_{k=0}^{\infty} C^k$. The error of the partial sum ending at C^k is at most $\|(I - C)^{-1}\|\cdot\|C\|^{k+1}$.
8. *Affine Iteration.* Suppose $C \in M_n(\mathbb{C})$ satisfies $\rho(C) < 1$, so $\|C\| < 1$ for some induced matrix norm. For any $d, x_0 \in \mathbb{C}^n$, the iteration $x_{k+1} = Cx_k + d$ for $k \geq 0$ will converge to the unique $x \in \mathbb{C}^n$ solving $(I - C)x = d$, and $\|x_k - x\| \leq \|C\|^k\|x_0 - x\|$.
9. *Richardson's Algorithm.* This iterative method solves $Ax = b$ by the update rule $x_{k+1} = x_k + (b - Ax_k)$. It will converge if $\rho(I - A) < 1$ (i.e., $\|I - A\| < 1$ for some induced matrix norm), in which case $\|x_k - x\| \leq \|I - A\|^k\|x_0 - x\|$. For instance, convergence is guaranteed if $\max_i \sum_j |(I - A)(i, j)| < 1$ or if $\max_j \sum_i |(I - A)(i, j)| < 1$.
10. *Jacobi's Algorithm.* This iterative method solves $Ax = b$ by the update rule $x_{k+1} = Q^{-1}(b - Rx_k)$ where Q is the diagonal part of A and $R = A - Q$. In detail,

$$x_{k+1}(i) = A(i, i)^{-1} \left[b(i) - \sum_{\substack{j=1 \\ j \neq i}}^n A(i, j)x_k(j) \right].$$

The method converges if $\rho(I - Q^{-1}A) < 1$, in which case $\|x_k - x\| \leq \|I - Q^{-1}A\|^k\|x_0 - x\|$. For instance, convergence is guaranteed if A is diagonally dominant ($|A(i, i)| < \sum_{j \neq i} |A(i, j)|$ for all $i \in [n]$).

11. *Gauss-Seidel Algorithm.* This iterative method solves $Ax = b$ by the update rule

$Qx_{k+1} = b - Rx_k$ where Q is the lower-triangular part of A and $R = A - Q$. In detail,

$$x_{k+1}(i) = A(i, i)^{-1} \left[b(i) - \sum_{j < i} A(i, j)x_{k+1}(j) - \sum_{j > i} A(i, j)x_k(j) \right].$$

The method converges if $\rho(I - Q^{-1}A) < 1$, in which case $\|x_k - x\| \leq \|I - Q^{-1}A\|^k \|x_0 - x\|$. For instance, convergence is guaranteed if A is diagonally dominant.

12. *Power Method.* The power method approximates the largest eigenvalue c of $A \in M_n(\mathbb{C})$ by starting with $x_0 \neq 0$, computing $x_{k+1} = Ax_k/\|Ax_k\|_\infty$ for $k \geq 0$, and estimating $c \approx (Ax_k)(i)/x_k(i)$ where i is chosen so $|x_k(i)|$ is not too close to zero. If A is diagonalizable with a unique eigenvalue c of magnitude $\rho(A)$, then for most choices of x_0 , the power method approximations will converge to c , and the x_k 's will approach associated eigenvectors (depending on k). (The requirement on x_0 is that the expansion of x_0 as a linear combination of eigenvectors of A must involve the eigenvector associated with c with nonzero coefficient.)
13. *Variations of the Power Method.* The *shifted* power method finds the eigenvalue of $A \in M_n(\mathbb{C})$ farthest from $b \in \mathbb{C}$ by applying the power method to $A - bI$ and adding b to the output. The *inverse* power method finds the eigenvalue of A closest to zero by applying the power method to A^{-1} and taking the reciprocal of the output. The *shifted inverse* power method finds the eigenvalue of A closest to b by using the power method to get the largest eigenvalue c of $(A - bI)^{-1}$, and returning $c^{-1} + b$.
14. *Deflation.* If we have found one eigenvalue c of $A \in M_n(\mathbb{C})$ and an associated eigenvector x , we can compute a matrix $B \in M_{n-1}(\mathbb{C})$ whose eigenvalues are the remaining eigenvalues of A as follows. Extend x to an orthonormal basis of \mathbb{C}^n by the Gram–Schmidt algorithm, and let U be the unitary matrix having these basis vectors as columns. Let B be the matrix U^*AU with the first row and column erased.

10.19 Exercises

1. (a) Given $v = (2, -4, -1, 0, 3)$, compute $\|v\|_\infty$, $\|v\|_1$, and $\|v\|_2$. (b) Given $v = (3 + 4i, -2 - i, 5i)$, compute $\|v\|_\infty$, $\|v\|_1$, and $\|v\|_2$.
2. (a) Sketch the set of unit vectors in \mathbb{R}^2 for each of the norms $\|\cdot\|_\infty$, $\|\cdot\|_1$, and $\|\cdot\|_2$. (b) Repeat (a) in \mathbb{R}^3 .
3. For each matrix A , execute Richardson's method (computing x_1, \dots, x_5 by hand) to try to solve $Ax = b$ using $b = [1 \ 2]^T$ and $x_0 = 0$. Compare x_5 to the exact solution x . Does the method appear to be converging? (a) $A = \begin{bmatrix} 1 & 1/2 \\ 0 & 1 \end{bmatrix}$; (b) $A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$; (c) $A = \begin{bmatrix} 1 & 0.1 \\ 0.1 & 1 \end{bmatrix}$; (d) $A = \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}$.
4. For each matrix A , predict whether Richardson's method will converge. If so, give an estimate on the error $\|x_k - x\|$ for an appropriate vector norm.

(a) $A = \begin{bmatrix} 1.3 & 0.2 & -0.3 \\ 0.1 & 0.9 & -0.2 \\ 0.2 & 0.2 & 1.2 \end{bmatrix}$; (b) $A = \begin{bmatrix} 1 & 1 & 2 \\ 0 & 1 & 2 \\ 2 & 1 & 1 \end{bmatrix}$; (c) $A \in M_n(\mathbb{R})$ given by $A(i,j) = 2^{i-j}$ for $i \leq j$, and $A(i,j) = 0$ for $i > j$; (d) $A = \begin{bmatrix} 3.5 & -4.5 \\ 1.8 & -2.2 \end{bmatrix}$.

5. Write a program in a computer algebra system to implement Richardson's method. Use the program to solve $Ax = b$, where

$$A = \begin{bmatrix} 1 & 0.1 & 0 & 0.2 & -0.1 & 0.1 \\ 0.1 & 1 & 0 & 0.1 & 0 & 0.2 \\ 0 & 0.2 & 1 & -0.1 & 0 & 0.1 \\ -0.3 & 0 & 0 & 1 & -0.1 & 0 \\ -0.1 & 0.1 & 0 & 0 & 1 & 0 \\ 0 & 0.2 & 0.1 & -0.1 & 0.3 & 1 \end{bmatrix} \text{ and } b = \begin{bmatrix} 2 \\ 1.3 \\ 0 \\ -1 \\ 0.6 \\ 1.7 \end{bmatrix}.$$

Give an error estimate for your answer.

6. For each matrix A and column vector b , execute Jacobi's method (starting with $x_0 = 0$ and computing x_1, \dots, x_5) to try to solve $Ax = b$. Compare x_5 to the exact solution x . Does the method appear to be converging? (a) $A = \begin{bmatrix} 2 & 1 \\ -1 & 2 \end{bmatrix}$, $b = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$; (b) $A = \begin{bmatrix} 0 & 5 \\ -2 & 1 \end{bmatrix}$, $b = \begin{bmatrix} 2 \\ -2 \end{bmatrix}$; (c) $A = \begin{bmatrix} 1 & 3 \\ 3 & 1 \end{bmatrix}$, $b = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$; (d) $A = \begin{bmatrix} 3 & 1 & 1 \\ 2 & 4 & -1 \\ 0 & -1 & 2 \end{bmatrix}$, $b = \begin{bmatrix} 4 \\ -1 \\ 0 \end{bmatrix}$.
7. For each matrix A , predict whether Jacobi's method will converge. If so, give an estimate on the error $\|x_k - x\|$ for an appropriate vector norm.
- (a) $A = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 3 & 5 \\ 3 & 4 & 5 \end{bmatrix}$; (b) $A = \begin{bmatrix} 8 & 4 & 2 \\ 3 & 9 & 1 \\ 1 & -2 & 4 \end{bmatrix}$; (c) $A \in M_n(\mathbb{R})$ given by $A(i,i) = c > 2$ for $i \in [n]$, $A(i,j) = 1$ for $|i - j| = 1$, and $A(i,j) = 0$ for other i, j ; (d) $A = \begin{bmatrix} 2 & -5 \\ 1 & 3 \end{bmatrix}$.
8. Write a program in a computer algebra system to implement Jacobi's method. Use the program to solve $Ax = b$, where

$$A = \begin{bmatrix} 3.2 & 1.1 & 0 & 0 & -1.2 & 0 \\ 1.3 & 4.1 & -1.2 & 1.5 & 0 & 0 \\ 0.8 & 0 & -1.7 & 0 & 0.5 & 0 \\ 0.7 & 2.1 & 0 & 5.3 & 0 & -1.3 \\ 0.3 & 0 & -0.6 & 1.1 & -2.4 & 0 \\ 1.4 & -0.8 & 0 & 0 & -1.0 & 4.9 \end{bmatrix} \text{ and } b = \begin{bmatrix} 4.1 \\ -2.2 \\ 1.3 \\ -0.8 \\ 0 \\ -3.2 \end{bmatrix}.$$

Give an error estimate for your answer.

9. Repeat Exercise 6 using the Gauss-Seidel algorithm.
10. Repeat Exercise 7 using the Gauss-Seidel algorithm.
11. Repeat Exercise 8 using the Gauss-Seidel algorithm.
12. For each matrix A , execute several iterations of the power method to approximate the largest eigenvalue of A and an associated eigenvector. Also find all eigenvalues

of A exactly and compare to the algorithm's output. (a) $A = \begin{bmatrix} 4 & -2 \\ -3 & -1 \end{bmatrix}$;

(b) $A = \begin{bmatrix} 3 & 1 & 0 \\ 1 & 3 & 1 \\ 0 & 1 & 3 \end{bmatrix}$; (c) $A = \begin{bmatrix} 1 & 2 & i \\ 2 & 1 & -1 \\ -i & -1 & 1 \end{bmatrix}$.

13. Repeat Exercise 12, but use the inverse power method to find the smallest eigenvalue of A and an associated eigenvector.
14. Repeat Exercise 12, but use the shifted power method to find the eigenvalue of A farthest from 4 and an associated eigenvector.
15. Use the inverse shifted power method to find the eigenvalue of $A = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & -1 \\ 1 & -1 & 1 & 1 \\ 1 & -1 & 1 & -1 \end{bmatrix}$ closest to -3 and an associated eigenvector.
16. For each matrix A , find the exact eigenvalues of A (with multiplicities). Try executing the power method on A and x_0 for a few iterations, and explain what goes wrong (and why) in each case. (a) $A = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$, $x_0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$;
- (b) $A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$, $x_0 = \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}$; (c) $A = \begin{bmatrix} 3 & 2 & 1 \\ 0 & 2 & 1 \\ 0 & 0 & 1 \end{bmatrix}$, $x_0 = \begin{bmatrix} 0 \\ -3 \\ 4 \end{bmatrix}$;
- (d) $A = \begin{bmatrix} -1/2 & \sqrt{3}/2 \\ -\sqrt{3}/2 & -1/2 \end{bmatrix}$, $x_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$.
17. Write a program in a computer algebra system to implement the power method for calculating eigenvalues. Test your program on the matrices in Exercises 5 and 8.
18. For each matrix A in Exercise 3, try to approximate A^{-1} using the geometric series formula (10.11). For which matrices does the series converge?
19. Use the geometric series formula (10.11) to estimate A^{-1} , where A is the matrix in Exercise 5.
20. (a) For fixed $b \in \mathbb{C}$ and $n \in \mathbb{N}^+$, define $A \in M_n(\mathbb{C})$ by setting $A(i,j) = 1$ for $i = j$, $A(i,j) = b$ for $j = i + 1$, and $A(i,j) = 0$ for all other i,j . Use (10.11) to find A^{-1} . (b) Generalize (a) by showing that if C is nilpotent ($C^k = 0$ for some $k \in \mathbb{N}^+$), then $I - C$ is invertible and $(I - C)^{-1}$ is given exactly by the finite sum $I + C + C^2 + \cdots + C^{k-1}$.
21. Consider the affine iteration $x_{k+1} = Cx_k + d$ for fixed $C \in M_n(\mathbb{C})$ and $d \in \mathbb{C}^n$.
 - (a) Express x_k as a sum of terms involving powers of C , x_0 , and d .
 - (b) Prove that if C is nilpotent with $C^r = 0$, then x_r is the exact solution to $(I - C)x = d$.
22. Suppose $C \in M_n(\mathbb{C})$ satisfies $\rho(C) > 1$ and $I - C$ is invertible.
 - (a) Show that for all $d \in \mathbb{C}^n$, there exists $x_0 \in \mathbb{C}^n$ such that the affine iteration (10.12) starting at x_0 will converge.
 - (b) Show there exist $d, x_0 \in \mathbb{C}^n$ such that the affine iteration (10.12) starting at x_0 will not converge.
23. Prove that the 1-norm on \mathbb{C}^n satisfies the axioms for a vector norm.
24. For $x = (x_1, x_2, x_3) \in \mathbb{R}^3$, define $\|x\| = |x_1| + |x_3|$. Which axioms for a vector norm are satisfied?
25. For $x = (x_1, x_2) \in \mathbb{R}^2$, define $\|x\| = (\sqrt{|x_1|} + \sqrt{|x_2|})^2$. Which axioms for a vector norm are satisfied?

26. *Cauchy–Schwarz Inequality.* For $x, y \in \mathbb{R}^n$, prove $|x \bullet y| \leq \|x\|_2 \|y\|_2$. [Hint: compute $(ax + by) \bullet (ax + by)$ for $a, b \in \mathbb{R}$. For $x \neq 0 \neq y$, take $a = \|x\|_2^{-1}$ and $b = \pm \|y\|_2^{-1}$.]
27. Prove that the 2-norm on \mathbb{R}^n satisfies the axioms for a vector norm. (For the triangle inequality, compute $\|x + y\|_2^2$.)
28. For $x = (x_1, x_2) \in \mathbb{R}^2$, define $\|x\| = |x_1|^2 + |x_2|^2$. Which axioms for a vector norm are satisfied?
29. Let $(V, \|\cdot\|_V)$ and $(W, \|\cdot\|_W)$ be normed vector spaces. (a) For $(v, w) \in V \times W$, define $\|(v, w)\| = \|v\|_V + \|w\|_W$. Prove this defines a vector norm on $V \times W$. (b) For $(v, w) \in V \times W$, define $\|(v, w)\| = \max(\|v\|_V, \|w\|_W)$. Prove this defines a vector norm on $V \times W$. (c) Explain why the 1-norm and sup norm on \mathbb{R}^n are special cases of the constructions in (a) and (b).
30. For each real p with $1 \leq p < \infty$, the p -norm on \mathbb{C}^n is defined by setting $\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$ for $x = (x_1, \dots, x_n) \in \mathbb{C}^n$. (a) Show $\|\cdot\|_p$ satisfies the first two axioms for a vector norm. (b) *Hölder's Inequality.* Assume $1 < p, q < \infty$ and $p^{-1} + q^{-1} = 1$. Prove: for all $x, y \in \mathbb{C}^n$, $\sum_{i=1}^n |x_i y_i| \leq \|x\|_p \|y\|_q$. (First prove it assuming $\|x\|_p = \|y\|_q = 1$, using Exercise 74 of Chapter 11.) (c) *Minkowski's Inequality.* Prove: for all $x, y \in \mathbb{C}^n$, $\|x + y\|_p \leq \|x\|_p + \|y\|_p$. (Start with $\|x + y\|_p^p$, observe that $|x_i + y_i|^p \leq (|x_i| + |y_i|) \cdot |x_i + y_i|^{p-1}$, and use (b).)
31. Prove: for $x \in \mathbb{R}^n$, $\lim_{p \rightarrow \infty} \|x\|_p = \|x\|_\infty$ (see Exercise 30 for the definition of $\|x\|_p$).
32. For fixed real numbers $p, r \geq 1$, find the minimum constant $C = C(p, r)$ such that for all $x \in \mathbb{C}^n$, $\|x\|_p \leq C \|x\|_r$. (One approach is to use Lagrange multipliers.)
33. Prove: for all x, y in a normed vector space, $\| \|x\| - \|y\| \| \leq \|x - y\|$.
34. (a) Fix an invertible $S \in M_n(\mathbb{C})$. Define $\|v\|_S = \|S^{-1}v\|_\infty$ for $v \in \mathbb{C}^n$. Prove this is a vector norm on \mathbb{C}^n . (b) Prove that the matrix norm induced by $\|\cdot\|_S$ satisfies $\|A\|_S = \|S^{-1}AS\|_\infty$ for all $A \in M_n(\mathbb{C})$.
35. Suppose $T : V \rightarrow W$ is a vector space isomorphism. (a) Given a vector norm $\|\cdot\|_V$ on V , show that setting $\|w\|_T = \|T^{-1}(w)\|_V$ for $w \in W$ defines a vector norm on W . (b) Prove that for the metrics associated with the norms defined in (a), $d(T(v), T(v')) = d(v, v')$ for all $v, v' \in V$. (c) Explain why the vector norm $\|\cdot\|_S$ in §10.13 is a special case of the construction in this exercise.
36. Fix $i \in [n]$. (a) Prove directly from the definitions that for all $x = (x_1, \dots, x_n) \in \mathbb{C}^n$, $|x_i| \leq \|x\|$ where $\|\cdot\|$ is the 1-norm, the 2-norm, or the sup norm. (b) Use a theorem from the text to show that for any vector norm on \mathbb{C}^n , there is a constant K with $|x_i| \leq K \|x\|$. (c) Let V be a finite-dimensional real or complex normed vector space with ordered basis $Z = (z_1, \dots, z_n)$. Each $v \in V$ has a unique expansion $v = p_1(v)z_1 + \dots + p_n(v)z_n$ for certain scalars $p_i(v)$. Prove there are constants K_i such that for all $v \in V$, $|p_i(v)| \leq K_i \|v\|$.
37. Consider the \mathbb{Q} -vector space $V = \{a + b\sqrt{2} : a, b \in \mathbb{Q}\}$. (a) Define $\|a + b\sqrt{2}\|$ to be the absolute value (in \mathbb{R}) of $a + b\sqrt{2}$. Show that the axioms for a vector norm (using rational scalars) hold. (b) Show that there does not exist a constant K such that for all $x = a + b\sqrt{2} \in V$, $|a| \leq K \|x\|$. (Contrast to Exercise 36(c).)
38. Let V be the infinite-dimensional real vector space of continuous functions $f : [0, 1] \rightarrow \mathbb{R}$. (a) Show that $\|f\|_\infty = \sup_{x \in [0, 1]} |f(x)|$ defines a vector norm on V . (b) Show that $\|f\|_I = \int_0^1 |f(x)| dx$ defines a vector norm on V .

39. Let V , $\|\cdot\|_\infty$, and $\|\cdot\|_I$ be defined as in Exercise 38. (a) Show there is no constant C such that for all $f \in V$, $\|f\|_\infty \leq C\|f\|_I$. (b) Is there a constant C such that for all $f \in V$, $\|f\|_I \leq C\|f\|_\infty$? Explain.
40. Let $n > 1$. (a) For $A \in M_n(\mathbb{C})$, define $\|A\| = \max_{i,j \in [n]} |A(i,j)|$. Prove that this is not a matrix norm. (b) For $A \in M_n(\mathbb{C})$, define $\|A\| = \sum_{i,j \in [n]} |A(i,j)|$. Is this a matrix norm? If so, is this an induced matrix norm?
41. *Frobenius Matrix Norm.* For $A \in M_n(\mathbb{C})$, define $\|A\|_F = \sqrt{\sum_{i,j \in [n]} |A(i,j)|^2}$.
 (a) Prove $\|\cdot\|_F$ is a matrix norm. (b) Is $\|\cdot\|_F$ induced by any vector norm on \mathbb{C}^n ? Explain. (c) Prove $\|A\|_F^2$ is the trace (sum of diagonal entries) of A^*A .
42. Let U be a unitary matrix. (a) Prove: for all $v \in \mathbb{C}^n$, $\|Uv\|_2 = \|v\|_2$. (b) Must (a) hold for the 1-norm or the sup norm on \mathbb{C}^n ? (c) Prove: for all $A \in M_n(\mathbb{C})$, $\|UA\|_2 = \|A\|_2 = \|AU\|_2$. (d) Does (c) hold for the matrix norms induced by the 1-norm or the sup norm on \mathbb{C}^n ? (e) Does (c) hold for the Frobenius matrix norm (Exercise 41)?
43. For any real or complex vector space V and $x, y \in V$, let $d(x, y) = 1$ if $x \neq y$, and $d(x, y) = 0$ if $x = y$. (a) Prove that (V, d) is a metric space. (b) Is there a vector norm on V such that d is the associated metric? Explain.
44. Suppose V is a real or complex vector space and d is a metric on V invariant under translations and respecting dilations. For $x \in V$, define $\|x\| = d(x, 0_V)$. Prove $\|\cdot\|$ is a vector norm whose associated metric is d .
45. Let (X, d_X) and (Y, d_Y) be metric spaces. Show $X \times Y$ is a metric space with metric given by $d((x_1, y_1), (x_2, y_2)) = d_X(x_1, x_2) + d_Y(y_1, y_2)$ for $x_1, x_2 \in X$ and $y_1, y_2 \in Y$.
46. Suppose (X, d) is a metric space and $(x_k : k \geq 0)$ is a sequence in X . (a) Suppose $x_k = x \in X$ for all $k \in \mathbb{N}$. Prove $\lim_{k \rightarrow \infty} x_k = x$. (b) Assume $x_k \rightarrow x \in X$. Prove: for all $j \in \mathbb{N}$, $\lim_{k \rightarrow \infty} x_{k+j} = x$. (c) Suppose $(y_k : k \geq 0)$ is a sequence such that $y_k = x_{j_k}$ for some indices $j_0 < j_1 < \dots < j_k < \dots$ (we say (y_k) is a *subsequence* of (x_k)). Prove: if $x_k \rightarrow x$, then $y_k \rightarrow x$. (d) Give an example to show that the converse of the result in (c) is false in general.
47. Suppose $(x_k : k \geq 0)$ is a sequence in a normed vector space V . (a) Prove: if $x_k \rightarrow x$ in V , then $\|x_k\| \rightarrow \|x\|$ in \mathbb{R} . (b) Show that if $\|x_k\| \rightarrow 0$ in \mathbb{R} , then $x_k \rightarrow 0_V$ in V . (c) Give an example to show that if $\|x_k\| \rightarrow r > 0$ in \mathbb{R} , then (x_k) need not converge in V .
48. Let (X, d) be a metric space with $x_k \in X$ for $k \geq 0$. We say (x_k) is a *Cauchy sequence* iff for all $\epsilon > 0$, there exists $k_0 \in \mathbb{N}^+$ such that for all $j, k \geq k_0$, $d(x_k, x_j) < \epsilon$. (a) Prove that every convergent sequence is a Cauchy sequence. (b) Give an example of a Cauchy sequence that does not converge. (Try the metric space $X = \mathbb{Q}$ with $d(x, y) = |x - y|$ for $x, y \in \mathbb{Q}$.) (c) Prove that if (x_k) is a Cauchy sequence and some subsequence $(y_k) = (x_{j_k})$ converges to x , then $x_k \rightarrow x$.
49. We say a metric space (X, d) is *complete* iff every Cauchy sequence in X converges to some point in X (see Exercise 48). (a) Given that \mathbb{R}^1 is complete, prove that \mathbb{R}^n is complete with the metric $d(x, y) = \|x - y\|_\infty$ for $x, y \in \mathbb{R}^n$. (b) For any vector norm $\|\cdot\|$ on \mathbb{R}^n , use (a) to prove that \mathbb{R}^n is complete with the metric $d(x, y) = \|x - y\|$ for $x, y \in \mathbb{R}^n$.
50. Suppose $c_k \rightarrow c$ in \mathbb{R} or \mathbb{C} and $v_k \rightarrow v$ in a normed vector space V . Prove $c_k v_k \rightarrow cv$.
51. Prove that $\|x\|_\infty \leq \|x\|_1 \leq n\|x\|_\infty$ for all $x \in \mathbb{R}^n$.

52. (a) Find a nonzero $x \in \mathbb{R}^n$ such that $\|x\|_1 = \|x\|_2 = \|x\|_\infty$. (b) Find all $y \in \mathbb{R}^n$ such that $\|y\|_2 = \sqrt{n}\|y\|_\infty$. (c) Find all $z \in \mathbb{R}^n$ such that $\|z\|_1 = n\|z\|_\infty$.
53. Let S be the set of all vector norms on a (possibly infinite-dimensional) vector space V . (a) Show that comparability of norms defines an equivalence relation on S . (b) How many equivalence classes are there when $V = \mathbb{R}^n$?
54. Let $\|\cdot\|$ be any vector norm on \mathbb{R}^n . (a) Show that the set of unit vectors $S = \{x \in \mathbb{R}^n : \|x\| = 1\}$ is a closed and bounded (hence compact) subset of the metric space (\mathbb{R}^n, d) , where $d(y, z) = \|y - z\|_2$ for $y, z \in \mathbb{R}^n$. (b) For the matrix norm $\|A\|$ induced by $\|\cdot\|$, show there exists $x \in S$ with $\|Ax\| = \|A\|$.
55. Let $\|\cdot\|$ be any matrix norm. (a) Prove: for all $A \in M_n(\mathbb{C})$ and all $k \in \mathbb{N}$, $\|A^k\| \leq \|A\|^k$. (b) If A is invertible, is there any relationship between $\|A^{-1}\|$ and $\|A\|^{-1}$?
56. Let $\|A\|$ be the matrix norm induced by a given vector norm $\|\cdot\|$ on \mathbb{C}^n . (a) Prove: $\|A\| = \sup\{\|Ax\|/\|x\| : x \in \mathbb{C}^n, x \neq 0\}$. (b) Prove: $\|A\| = \inf\{C \in \mathbb{R} : \|Ax\| \leq C\|x\| \text{ for all } x \in \mathbb{C}^n\}$.
57. For each matrix A in Exercise 4, compute $\|A\|_\infty$, $\|A\|_1$, $\|A\|_2$, and $\|A\|_F$ (see Exercise 41). Take $n = 4$ when computing $\|A\|_2$ in part (c).
58. (a) For $A \in M_n(\mathbb{C})$, show that $\|A\|_\infty = \|(|A_{[1]}|_1, \dots, |A_{[n]}|_1)\|_\infty$. (b) State a similar formula for $\|A\|_1$.
59. Let $S = \{u \in \mathbb{R}^n : \|u\|_2 = 1\}$. For each matrix A below, define $f : S \rightarrow \mathbb{R}$ by $f(u) = \|Au\|_2^2$ for $u \in S$. Find the maximum value of f on S and hence compute $\|A\|_2$. (Use Lagrange multipliers for (b) and (c).) (a) $A = \begin{bmatrix} 3 & 1 \\ -1 & 3 \end{bmatrix}$; (b) $A = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$; (c) $A = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \end{bmatrix}$.
60. Prove: for all $A \in M_n(\mathbb{C})$, $\|A\|_2 = \sqrt{\rho(A^*A)}$. [Hint: A^*A is positive semidefinite.]
61. Use the theorem in §7.7 to show that given any $\epsilon > 0$ and any $B \in M_n(\mathbb{C})$, there is an invertible $S \in M_n(\mathbb{C})$ such that $S^{-1}BS$ is upper-triangular and $\|S^{-1}BS\|_\infty \leq \rho(B) + \epsilon$. [Hint: If T is triangular and D is diagonal with $D(i, i) = \delta^i$ for some $\delta > 0$, compare the entries of T and $D^{-1}TD$.]
62. Assume the setup in §10.15. (a) Suppose the largest eigenvalue c has several associated eigenvectors, say $c = c_1 = \dots = c_m$ and $|c_m| > |c_{m+1}|$. Prove the power method still converges. (b) Suppose instead that $|c_1| = \dots = |c_m| > |c_{m+1}|$ with $m > 1$. Must the power method converge? (c) Suppose we happen to pick x_0 with $d_1 = 0$. Assuming $|c_2| > |c_3|$, what will the power method do?
63. With the setup in §10.15, prove that $x_k - b'_k z_1 \rightarrow 0$ as $k \rightarrow \infty$.
64. Execute the deflation algorithm on the matrix $A = \begin{bmatrix} 1 & 4 & -1 \\ 0 & 3 & 2 \\ 1 & -1 & -3 \end{bmatrix}$, which has a known eigenvalue $c = -1$ and associated eigenvector $x = [3 \ -1 \ 2]^T$.
65. Let $A = \begin{bmatrix} 0 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$. (a) Find the largest eigenvalue of A and an associated eigenvector by the power method. (b) Use deflation to recover the other two eigenvalues of A .

66. Let $A = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$. (a) By inspection, find an eigenvector of A associated with the eigenvalue 1. (b) Use deflation to find $B \in M_3(\mathbb{C})$ such that $\chi_B(t) = (t - 1)\chi_A(t)$.
67. (a) Give an example of a matrix A and a scalar c such that Richardson's algorithm fails to converge when applied to $Ax = b$, but does converge when applied to $(cA)x = cb$. (b) For any $A \in M_n(\mathbb{C})$ and nonzero $c \in \mathbb{C}$, how are the eigenvalues of $I - A$ related to the eigenvalues of $I - cA$? (c) Find a condition on the eigenvalues of A guaranteeing the existence of $c \in \mathbb{C}$ such that Richardson's algorithm will converge for the system $(cA)x = cb$. Describe how to choose c if we know all the eigenvalues of A .
68. Let $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$, $Q_1 = \begin{bmatrix} a & 0 \\ 0 & d \end{bmatrix}$, and $Q_2 = \begin{bmatrix} a & 0 \\ c & d \end{bmatrix}$ where $a, b, c, d \in \mathbb{C}$ with $a \neq 0 \neq d$. (a) Find the eigenvalues of $I - Q_1^{-1}A$ and $I - Q_2^{-1}A$. (b) Deduce sufficient conditions for the convergence of the Jacobi method and the Gauss-Seidel method on the system $Ax = y$.
69. Give a specific example of a matrix $A \in M_3(\mathbb{R})$ and $b \in \mathbb{R}^3$ such that the Jacobi method converges when applied to $Ax = b$, but the Gauss-Seidel method does not converge for this system.
70. Give a specific example of a matrix $A \in M_3(\mathbb{R})$ and $b \in \mathbb{R}^3$ such that the Gauss-Seidel method converges when applied to $Ax = b$, but the Jacobi method does not converge for this system.
71. Discuss how the shifted power method might be used to find the largest eigenvalue of $A \in M_n(\mathbb{C})$ in each of these cases: (a) A is not invertible; (b) two eigenvalues of A both have magnitude $\rho(A)$.
72. Suppose $A \in M_n(\mathbb{C})$ is lower-triangular with nonzero entries on the diagonal. (a) What happens if we use the Gauss-Seidel algorithm to solve $Ax = b$? (b) Must the Jacobi method converge for this system? Explain.
73. Define $\|A\| = \rho(A)$ for $A \in M_n(\mathbb{C})$. Which axioms for a matrix norm are satisfied?

Part IV

The Interplay of Geometry and Linear Algebra

This page intentionally left blank

11

Affine Geometry and Convexity

A central concept in linear algebra is the idea of a subspace of a vector space. Subspaces provide an algebraic abstraction of “uncurved” geometric objects such as lines and planes through the origin in three-dimensional space. However, there are many other uncurved geometric figures that are not subspaces, including lines and planes not passing through the origin, individual points, line segments, triangles, quadrilaterals, polygons in the plane, tetrahedra, cubes, and solid polyhedra. The concepts of *affine sets* and *convex sets* provide an algebraic setting for studying such figures and their higher-dimensional analogues.

Geometrically, an affine subset of a vector space is a subspace that has been “translated” away from the origin by adding some fixed vector. Affine sets can be described in other ways as well, e.g., as intersections of hyperplanes, as solution sets of systems of linear equations, or as the set of “affine combinations” of a given set of vectors. A notion of affine independence, analogous to the familiar idea of linear independence, provides affine versions of bases, dimension, and coordinate systems for affine sets. Affine maps preserve the affine structure of affine sets, just as linear maps preserve the linear structure of vector spaces and subspaces.

A convex set in \mathbb{R}^n contains the line segment joining any two of its points. Given a set S of points in \mathbb{R}^n , there is a smallest convex set containing S , called the *convex hull* of S . One can build up the convex hull from the generating set S , by taking so-called “convex combinations” of the generators, or one can obtain the convex hull by intersecting all convex sets that contain S . This fact is one instance of a family of theorems asserting the equivalence of a “generative” description of a class of sets and an “intersectional” description of that same class. We shall encounter several such theorems in this chapter, including the fundamental result that convex hulls of finite sets coincide with bounded intersections of finitely many closed half-spaces.

The chapter concludes with a short discussion of convex real-valued functions, which are closely related to convex sets. These functions play a prominent role in analysis, linear algebra, and optimization theory (notably linear programming). We prove Jensen’s inequality for convex functions and describe how to test convexity by examination of the first and second derivatives of a function.

11.1 Linear Subspaces

Before beginning our study of affine geometry, we want to discuss some facts about subspaces of vector spaces that will be helpful in the development of the affine theory. Let F be a field, let n be a positive integer, and let V be an n -dimensional vector space over F . We know V is isomorphic to the vector space F^n consisting of n -tuples $v = (v_1, \dots, v_n)$ with

each $v_i \in F$ (see §6.3). When convenient, we identify an n -tuple v with the $n \times 1$ column vector $\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$.

A *subspace* of V is a subset W of V satisfying these three closure conditions:

- (i) $0_V \in W$ (closure under identity);
- (ii) for all $v, w \in W$, $v + w \in W$ (closure under vector sum);
- (iii) for all $v \in W$ and $c \in F$, $cv \in W$ (closure under scalar multiplication).

In this chapter, we will often refer to subspaces of V as *linear subspaces* to emphasize the distinction between subspaces and affine sets (defined later). Given a subspace W , any $k \in \mathbb{N}$, vectors $v_1, \dots, v_k \in W$, and scalars $c_1, \dots, c_k \in F$, it follows by induction on k that $c_1v_1 + \dots + c_kv_k$ must belong to W . The expression $c_1v_1 + \dots + c_kv_k$ is called a *linear combination* of v_1, \dots, v_k , so we see that subspaces are closed under taking linear combinations of their elements. Conversely, any subset of V that is closed under linear combinations must be a linear subspace of V . (When proving the converse, one must keep in mind the convention that a linear combination with $k = 0$ terms is interpreted as 0_V .)

Geometrically, the condition that $0_V \in W$ means that all linear subspaces are required to “pass through the origin.” Every subspace W of V is itself an F -vector space, and the dimension of W is between 0 and n . One-dimensional subspaces of V are called *lines through the origin*; two-dimensional subspaces of V are called *planes through the origin*; and $(n - 1)$ -dimensional subspaces of the n -dimensional space V are called *linear hyperplanes*. Lines and planes in \mathbb{R}^3 that do not pass through $(0, 0, 0)$ are not linear subspaces; they are examples of the affine sets to be studied later.

11.2 Examples of Linear Subspaces

We now discuss four ways of constructing examples of linear subspaces, each of which leads to a different description of the points in the subspace. First, suppose that V is F^n (viewed as a set of column vectors), and $A \in M_{m,n}(F)$ is an $m \times n$ matrix with entries in F . The *null space* of A is $W = \{x \in V : Ax = 0\}$. We check that W is a subspace of V as follows: W is a subset of V by definition; $A0 = 0$, so (i) holds; for all $v, w \in W$, $Av = 0 = Aw$ implies $A(v + w) = Av + Aw = 0 + 0 = 0$, so (ii) holds; for all $v \in W$ and $c \in F$, $Av = 0$ implies $A(cv) = c(Av) = c0 = 0$, so (iii) holds. We can think of the null space of A more concretely as the *solution set of the system of m homogeneous linear equations in n unknowns*

$$\left\{ \begin{array}{rcl} A(1,1)x_1 + A(1,2)x_2 + \cdots + A(1,n)x_n & = & 0 \\ A(2,1)x_1 + A(2,2)x_2 + \cdots + A(2,n)x_n & = & 0 \\ \cdots & & \cdots \\ A(m,1)x_1 + A(m,2)x_2 + \cdots + A(m,n)x_n & = & 0. \end{array} \right. \quad (11.1)$$

The equations are called *homogeneous* because their right-hand sides are zero.

Second, suppose that V is F^m (viewed as column vectors), and $A \in M_{m,n}(F)$ is a given matrix. The *range* of A is $W = \{v \in V : \exists x \in F^n, Ax = v\}$. It is routine to verify that W is a subspace of F^m . We can give a more abstract version of these first two constructions by considering an arbitrary F -linear map $T : V \rightarrow Z$ between two F -vector spaces V and

Z . The *kernel* of T , defined to be $\ker(T) = \{v \in V : T(v) = 0_Z\}$, is a subspace of V . The *image* of T , defined to be $\text{img}(T) = \{z \in Z : \exists v \in V, T(v) = z\}$, is a subspace of Z . To recover the earlier examples involving the matrix A , take $T : F^n \rightarrow F^m$ to be $T(v) = Av$ for $v \in F^n$. This is a linear map whose kernel is the null space of A and whose image is the range of A .

Third, suppose $S = \{v_1, \dots, v_k\}$ is a set of k vectors in V . Let W consist of all linear combinations of elements of S , so that a vector w is in W iff there exist scalars $c_1, \dots, c_k \in F$ with $w = c_1v_1 + \dots + c_kv_k$. It is routine to check that W satisfies the closure conditions (i), (ii), and (iii), so that W is a subspace of V . We call W the *subspace spanned by S* or the *linear span of S* , writing $W = \text{Sp}(S)$ or $W = \text{Sp}_F(S)$. Observe that $S \subseteq W$, because each v_i can be written in the form $0v_1 + \dots + 1v_i + \dots + 0v_n$. For infinite $S \subseteq V$, let $\text{Sp}_F(S)$ consist of all linear combinations of finitely many elements of S . Again, one sees that this is a linear subspace of V containing S . For all linear subspaces Z of V , if $S \subseteq Z$ then $\text{Sp}_F(S) \subseteq Z$ since Z is closed under linear combinations. So, $\text{Sp}_F(S)$ is the smallest linear subspace of V containing S . The fact that the range of $A \in M_{m,n}(F)$ is a subspace is a special case of the spanning construction. This follows from the fact that the elements Ax in the range of A are precisely the linear combinations $x_1A^{[1]} + x_2A^{[2]} + \dots + x_nA^{[n]}$ of the n columns of A (see §4.7).

Fourth, suppose $\{W_i : i \in I\}$ is any collection of linear subspaces of V . The intersection $W = \bigcap_{i \in I} W_i$ is readily seen to be a subspace of V . For example, *any intersection of linear hyperplanes in V is a subspace of V* . We claim that the null space of a matrix $A \in M_{m,n}(F)$ can be viewed as a special case of this construction. To see why, note that the solution set of the system of equations (11.1) is the intersection of the solution sets of the m individual equations considered separately. Consider the solution set Z of a particular equation $c_1x_1 + \dots + c_nx_n = 0$, where every $c_i \in F$. If all c_i 's are zero, then Z is all of F^n , and we can discard this equation from the system without changing the final solution set of the full system. If some $c_i \neq 0$, we solve the equation for x_i to get

$$x_i = c_i^{-1}(-c_1x_1 - \dots - c_{i-1}x_{i-1} - c_{i+1}x_{i+1} - \dots - c_nx_n).$$

We obtain all solutions by choosing any scalar value for each x_j with $j \neq i$ and then using the displayed relation to compute x_i . For example, if $k \neq i$ is fixed and we choose $x_k = 1$ and $x_j = 0$ for all $j \neq k, i$, then $x_i = -c_i^{-1}c_k$. From this description, one can verify that the set $\{e_k - c_i^{-1}c_k e_i : 1 \leq k \leq n, k \neq i\}$ is a basis for Z of size $n - 1$, where e_i and e_k are standard basis vectors in F^n . This means that Z is a linear hyperplane in F^n . The null space of A is the intersection of the linear hyperplanes associated with each nonzero row of A . (If A is a zero matrix, we use the convention that the intersection of an empty collection of hyperplanes is all of F^n .)

11.3 Characterizations of Linear Subspaces

We claim that every linear subspace of V can be described in one of the ways discussed in the last section. More specifically, we show that: (a) every k -dimensional subspace W of V is the linear span of a set of k linearly independent vectors; (b) every $(n - m)$ -dimensional subspace W of F^n is the solution set of a system of m homogeneous nonzero linear equations $Ax = 0$ for some $A \in M_{m,n}(F)$ (i.e., W is the null space of A); (c) every $(n - m)$ -dimensional subspace W of V is an intersection of m linear hyperplanes in V .

Part (a) follows from the fact that the subspace W is itself a vector space, so it has a basis S of some size k (see §1.8 and Chapter 16). We must have $k \leq n$, since no linearly

independent set in V has size larger than $n = \dim(V)$. Since S spans W , W consists of all linear combinations of elements of S . In fact, the linear independence of S ensures that each $w \in W$ can be written in *exactly one* way as $c_1v_1 + \cdots + c_kv_k$ with $c_i \in F$ and $v_i \in S$.

For (b), let W be an $(n - m)$ -dimensional subspace of F^n with ordered basis $B = (v_{m+1}, v_{m+2}, \dots, v_n)$. By adding appropriate vectors to the beginning of the list B , we can extend B to an ordered basis $C = (v_1, \dots, v_m, v_{m+1}, \dots, v_n)$ for F^n . Define $T : F^n \rightarrow F^m$ by setting

$$T(c_1v_1 + \cdots + c_mv_m + c_{m+1}v_{m+1} + \cdots + c_nv_n) = \begin{bmatrix} c_1 \\ \vdots \\ c_m \end{bmatrix} \quad (c_1, \dots, c_n \in F).$$

Let $A \in M_{m,n}(F)$ be the matrix of T with respect to the standard ordered bases of F^n and F^m , so $T(x) = Ax$ for all $x \in F^n$ (see Chapter 6). On one hand, the kernel of T consists precisely of the linear combinations of vectors in B , which is the given subspace W . On the other hand, since $T(x) = Ax$, we see that the kernel of T is exactly the null space of A , which in turn is the solution set of the system of equations $Ax = 0$. Since the image of T (the range of A) is all of F^m , it is evident that none of the rows of A can be zero. So we have expressed W as the solution set of m homogeneous nonzero linear equations. It follows (as shown in the last section) that W is the intersection of the m linear hyperplanes determined by each of the m rows of A .

We have now proved the special case of (c) where W is a subspace of F^n . For a general n -dimensional vector space V with $(n - m)$ -dimensional subspace W , we know there is a vector space isomorphism $g : V \rightarrow F^n$. (For example, as seen in Chapter 6, we can obtain g by selecting an ordered basis of V and mapping each $v \in V$ to the coordinates of v relative to this basis.) Applying the isomorphism g to W gives a subspace $W' = g[W]$, which is an $(n - m)$ -dimensional subspace of F^n . So we can find m linear hyperplanes H'_1, \dots, H'_m in F^n with $W' = H'_1 \cap \cdots \cap H'_m$. Let $H_i = g^{-1}[H'_i]$ for $1 \leq i \leq m$, which is an $(n - 1)$ -dimensional subspace (linear hyperplane) in V since g^{-1} is an isomorphism. Then

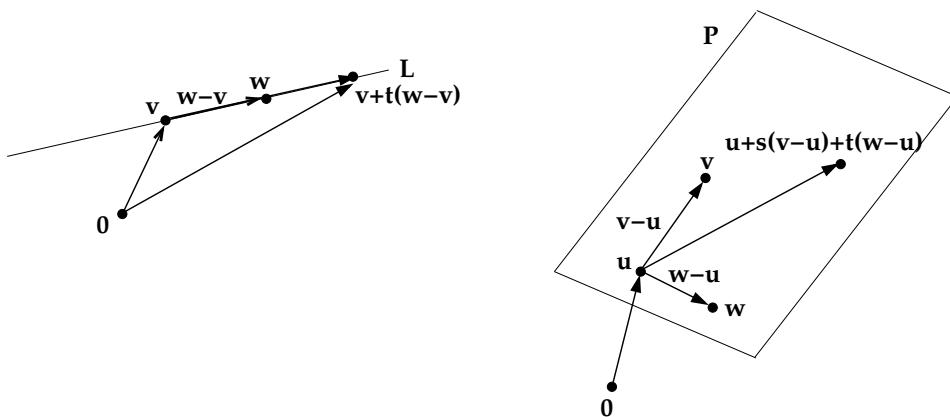
$$W = g^{-1}[W'] = g^{-1}[H'_1 \cap \cdots \cap H'_m] = g^{-1}[H'_1] \cap \cdots \cap g^{-1}[H'_m] = H_1 \cap \cdots \cap H_m,$$

so W is an intersection of m linear hyperplanes in V .

11.4 Affine Combinations and Affine Sets

We continue to assume that V is an n -dimensional vector space over a field F . We have seen that the linear subspaces of V are precisely the subsets of V closed under linear combinations. Now, define an *affine combination* of $v_1, \dots, v_k \in V$ to be a linear combination $c_1v_1 + \cdots + c_kv_k$ where each $c_i \in F$ and $c_1 + c_2 + \cdots + c_k = 1_F$. An *affine set* in V is a subset W of V that is closed under affine combinations. In other words, $W \subseteq V$ is affine iff for all $k \in \mathbb{N}^+$, $w_1, \dots, w_k \in W$ and $c_1, \dots, c_k \in F$, if $c_1 + \cdots + c_k = 1$ then $c_1w_1 + \cdots + c_kw_k \in W$. (For the field $F = \mathbb{R}$, one can prove that a set W is affine iff for all $w, z \in W$ and $c \in \mathbb{R}$, $cw + (1 - c)z \in W$, which corresponds to taking $k = 2$ above. But this is not true for all fields; see Exercise 19.)

To motivate the definition of affine combinations, we recall the vector equations of lines and planes in \mathbb{R}^3 that do not necessarily pass through the origin. Consider the line L in \mathbb{R}^3 through the distinct points $v, w \in \mathbb{R}^3$. Geometrically, we can reach every point on this line

**FIGURE 11.1**Vector Description of Lines and Planes in \mathbb{R}^3 .

by starting at the origin, traveling to the tip of the vector v (viewed as an arrow starting at 0), and then following some scalar multiple $t(w - v)$ of the “direction vector” $w - v$ for the line (where $t \in \mathbb{R}$). See Figure 11.1. Algebraically, the line consists of all vectors $v + t(w - v)$ as t ranges over \mathbb{R} . Restating this,

$$L = \{tw + (1 - t)v : t \in \mathbb{R}\} = \{tw + sv : t, s \in \mathbb{R} \text{ and } t + s = 1\}.$$

In other words, L consists of all affine combinations of v and w . One can check that L is an affine set in \mathbb{R}^3 (this is a special case of a general result proved in §11.6).

Similarly, consider the plane P in \mathbb{R}^3 through three non-collinear points $u, v, w \in \mathbb{R}^3$. As in the case of the line, we reach arbitrary points in P from the origin by going to the tip of u , then traveling within the plane some amount in the direction $v - u$ and some other amount in the direction $w - u$; see Figure 11.1. Algebraically,

$$\begin{aligned} P &= \{u + s(v - u) + t(w - u) : s, t \in \mathbb{R}\} = \{(1 - s - t)u + sv + tw : s, t \in \mathbb{R}\} \\ &= \{ru + sv + tw : r, s, t \in \mathbb{R}, r + s + t = 1\}. \end{aligned}$$

So P consists of all affine combinations of u, v, w , and one checks that P is an affine set in \mathbb{R}^3 .

A plane in \mathbb{R}^3 can also be specified as the solution set of one linear (possibly non-homogeneous) equation in three unknowns. A line in \mathbb{R}^3 is the solution set of two linear equations in three unknowns. More generally, the reader may check that for any matrix $A \in M_{m,n}(F)$ and $b \in F^m$, the set of all solutions $x \in F^n$ to the non-homogeneous system of linear equations $Ax = b$ is an affine subset of F^n . We shall see later (§11.9) that all affine subsets of F^n have this form.

11.5 Affine Sets and Linear Subspaces

We can use linear subspaces of V to build more examples of affine sets in V . Let W be a fixed linear subspace of V . Since W is closed under all linear combinations of its elements, which

include affine combinations as a special case, W is an affine subset of V . More generally, for any subset S of V and any vector $u \in V$, define the *translate of S by u* to be the set $u + S = \{u + x : x \in S\}$. Let us prove that each translate $u + W$ is an affine set. Fix $v_1, \dots, v_k \in u + W$ and $c_1, \dots, c_k \in F$ with $c_1 + \dots + c_k = 1$; we must check that $c_1v_1 + \dots + c_kv_k \in u + W$. Write each v_i as $u + w_i$ for some $w_i \in W$; then

$$\sum_{i=1}^k c_i v_i = \sum_{i=1}^k c_i(u + w_i) = \sum_{i=1}^k c_i u + \sum_{i=1}^k c_i w_i = u \sum_{i=1}^k c_i + \sum_{i=1}^k c_i w_i.$$

We know $\sum_{i=1}^k c_i = 1$ and $\sum_{i=1}^k c_i w_i \in W$ (because W is a linear subspace), so our calculation shows that $\sum_{i=1}^k c_i v_i$ is in $u + W$, as needed. In fact, exactly the same calculation proves that any translate $u + X$ of an affine set X is also an affine set.

Taking W to be the subspace V , we see that the whole space V is an affine set. Taking W to be the subspace $\{0\}$, we see that every one-point set $u + W = \{u\}$ is an affine set. By definition, an *affine line* is a translate $u + W$ where $\dim(W) = 1$; an *affine plane* is a translate $u + W$ where $\dim(W) = 2$; and an *affine hyperplane* is a translate $u + W$ where $\dim(W) = n - 1$ (i.e., W is a linear hyperplane).

Let us prove that *for all $X \subseteq V$, X is a linear subspace iff X is an affine set and $0_V \in X$.* We already know the forward direction is true. Conversely, suppose X is an affine set containing zero. Given $v, w \in X$, $v + w$ is the affine combination $-1 \cdot 0_V + 1v + 1w$ of elements of X , so $v + w \in X$. Given $v \in X$ and $c \in F$, cv is the affine combination $(1 - c) \cdot 0_V + cv$ of elements of X , so $cv \in X$. We also know $0_V \in X$, so X is a linear subspace.

Examining the definition closely, one sees that the empty set \emptyset is an affine subset of V . We now prove that *for every nonempty affine set X in V , there exists a unique linear subspace W in V such that $X = u + W$ for some $u \in V$.* The vector u here is not unique; in fact, we will see that $X = u + W$ and $X = -u + W$ for any choice of $u \in X$. We call W the *direction subspace of X* . This theorem provides one geometric characterization of the nonempty affine sets in V : they are precisely the translates of the linear subspaces of V .

To prove the theorem, fix a nonempty affine set X in V , and fix $u \in X$. Define $W = -u + X$. On one hand, W is an affine set, being a translate of the affine set X . On the other hand, $u \in X$ implies $0 = -u + u \in W$. By our earlier theorem, W is a linear subspace, and evidently $u + W = u + (-u + X) = X$. We now prove uniqueness of the linear subspace W . Say $X = u + W = v + W'$ for some $u, v \in V$ and linear subspaces W and W' . Since $0 \in W$ and $0 \in W'$, we have $u, v \in X$. Then $v = u + w$ for some $w \in W$, hence $W' = -v + X = (-v + u) + W = -w + W = W$. (The last step uses the fact that $x + W = W$ for all x in a subspace W .)

Since the direction subspace of an affine set X is uniquely determined by X , we can unambiguously define the *affine dimension* of X , denoted $\dim(X)$, to be the dimension (as a vector space) of its direction subspace. For instance, points, affine lines, affine planes, and affine hyperplanes have respective affine dimensions 0, 1, 2, and $n - 1$. The affine dimension of \emptyset is undefined.

11.6 Affine Span of a Set

Let $S = \{v_1, \dots, v_k\}$ be a finite subset of V . We know that the smallest linear subspace of V containing S is the linear span $\text{Sp}(S)$ consisting of all linear combinations of v_1, \dots, v_k .

By analogy, define the *affine span* or *affine hull* of S , denoted $\text{aff}(S)$, to be the set of all affine combinations of elements of S . In symbols,

$$\text{aff}(\{v_1, \dots, v_k\}) = \left\{ c_1 v_1 + \dots + c_k v_k : c_i \in F, \sum_{i=1}^k c_i = 1 \right\}.$$

Similarly, for an infinite subset S of V , let $\text{aff}(S)$ be the set of all affine combinations of finitely many vectors drawn from S . For example, in $V = \mathbb{R}^3$, our discussion of Figure 11.1 shows that $\text{aff}(\{v, w\})$ is the line through v and w , whereas $\text{aff}(\{u, v, w\})$ is the plane through u , v , and w .

We assert that $\text{aff}(S)$ is an affine set containing S , and it is the smallest affine set that contains S . To check that $\text{aff}(S)$ is affine, fix $w_1, \dots, w_s \in \text{aff}(S)$ and $c_1, \dots, c_s \in F$ with $\sum_{j=1}^s c_j = 1$. There are finitely many elements $v_1, \dots, v_k \in S$ such that each w_j is some affine combination of these elements, say $w_j = \sum_{i=1}^k a_{ij} v_i$ with $a_{ij} \in F$ and $\sum_{i=1}^k a_{ij} = 1$ for $1 \leq j \leq s$. Compute

$$w = \sum_{j=1}^s c_j w_j = \sum_{j=1}^s c_j \left(\sum_{i=1}^k a_{ij} v_i \right) = \sum_{j=1}^s \sum_{i=1}^k c_j a_{ij} v_i = \sum_{i=1}^k \left(\sum_{j=1}^s c_j a_{ij} \right) v_i.$$

So w is a linear combination of the v_i 's, and the coefficients in this combination have sum

$$\sum_{i=1}^k \sum_{j=1}^s c_j a_{ij} = \sum_{j=1}^s c_j \sum_{i=1}^k a_{ij} = \sum_{j=1}^s c_j \cdot 1 = 1.$$

Hence w is an affine combination of v_1, \dots, v_k , which proves that $w \in \text{aff}(S)$. For each $x \in S$, $x = 1x$ is an affine combination of elements of S , so $S \subseteq \text{aff}(S)$. Finally, suppose T is any other affine subset of V with $S \subseteq T$. Since T is closed under affine combinations of any of its elements, affine combinations of elements drawn from the subset S must lie in T . So $\text{aff}(S) \subseteq T$. This inclusion holds for all affine sets T containing S , one of which is $\text{aff}(S)$. It follows that $\text{aff}(S)$ is the intersection of all affine subsets of V that contain S .

11.7 Affine Independence

In three-dimensional geometry, any two distinct points determine a line, and any three non-collinear points determine a plane. We would like similar descriptions of general affine sets A of the form $A = \text{aff}(S)$, where the affine spanning set S is as small as possible. The examples of lines and planes suggest that the size of S should be one more than the affine dimension of A . To make this precise, we introduce the notion of affinely independent sets.

Recall that a list of vectors $L = (v_1, \dots, v_k)$ in V is *linearly independent* iff the only linear combination of the v_i 's that gives zero is the combination with all zero coefficients; i.e., for all $c_1, \dots, c_k \in F$, if $c_1 v_1 + \dots + c_k v_k = 0_V$ then every $c_i = 0_F$. We now define L to be *affinely independent* iff for all $c_1, \dots, c_k \in F$ such that $c_1 + c_2 + \dots + c_k = 0_F$, $c_1 v_1 + \dots + c_k v_k = 0_V$ implies every $c_i = 0_F$. The list L is *affinely dependent* iff there exist $c_1, \dots, c_k \in F$ with some $c_i \neq 0$, $\sum_{j=1}^k c_j = 0$, and $\sum_{j=1}^k c_j v_j = 0$. A set S of vectors (finite or not) is affinely independent iff every finite list of distinct vectors in S is affinely independent.

The following result relates affine independence to linear independence: *the list*

$L = (v_0, v_1, v_2, \dots, v_k)$ is affinely independent iff the list $L' = (v_1 - v_0, v_2 - v_0, \dots, v_k - v_0)$ is linearly independent. We prove the contrapositive in both directions. Assuming L is affinely dependent, we have scalars $c_0, \dots, c_k \in F$ summing to zero and not all zero, such that $\sum_{i=0}^k c_i v_i = 0$. We can rewrite the given combination of the v_i 's as

$$c_1(v_1 - v_0) + c_2(v_2 - v_0) + \cdots + c_k(v_k - v_0) + (c_0 + c_1 + c_2 + \cdots + c_k)v_0 = 0.$$

But the sum of all the c_i 's is zero, so we have expressed zero as a linear combination of the entries of L' with coefficients c_1, \dots, c_k . These coefficients cannot all be zero, since otherwise $\sum_{i=0}^k c_i = 0$ would give $c_0 = 0$ as well. This proves the linear dependence of L' . Conversely, assuming L' is linearly dependent, we have scalars $d_1, \dots, d_k \in F$ (not all zero) with $d_1(v_1 - v_0) + \cdots + d_k(v_k - v_0) = 0$. Defining $d_0 = -d_1 - d_2 - \cdots - d_k$, we then have $d_0 v_0 + d_1 v_1 + \cdots + d_k v_k = 0$ where not all d_i 's are zero, but the sum of all d_i 's is zero. So L is affinely dependent.

When using this result to detect affine independence of a finite set of vectors, we can list the elements of the set in any order. So, given a finite set $S \subseteq V$, we can test the affine independence of S by picking any convenient $x \in S$ and checking if the set $T = \{y - x : y \in S, y \neq x\}$ is linearly independent. For example, if S is any set of linearly independent vectors in V , then $S \cup \{0_V\}$ is affinely independent.

11.8 Affine Bases and Barycentric Coordinates

Let X be a nonempty affine subset of V . A list $L = (v_0, \dots, v_k)$ is an *ordered affine basis* of X iff L is affinely independent and $X = \text{aff}(\{v_0, \dots, v_k\})$. Similarly, a set S is an *affine basis* of X iff S is affinely independent and $X = \text{aff}(S)$.

We can use the direction subspace of X to find an affine basis for X . Specifically, write the given affine set X as $X = u + W$, where $u \in X$ and W is a linear subspace of V . Suppose $\dim(W) = k$, and let (w_1, \dots, w_k) be an ordered (linear) basis of W . We claim $L = (u, u + w_1, \dots, u + w_k)$ is an ordered affine basis of X . Using the criterion from §11.7, we see that L is affinely independent, since subtracting u from every other vector on the list produces the linearly independent list (w_1, \dots, w_k) . Next, $u \in X$ and each $u + w_i \in X$, so the set $S = \{u, u + w_1, \dots, u + w_k\}$ is contained in the affine set X . Hence $\text{aff}(S) \subseteq X$. To establish the reverse inclusion, fix any $z \in X$. Then $z - u \in W$, so we can write $z - u = c_1 w_1 + \cdots + c_k w_k$ for some $c_i \in F$. Hence, z itself can be written

$$z = u + \sum_{i=1}^k c_i w_i = \left(1 - \sum_{i=1}^k c_i\right)u + \sum_{i=1}^k c_i(u + w_i),$$

which is an affine combination of the elements of L . We now see that X , which has affine dimension $k = \dim(W)$, has an ordered affine basis consisting of $k + 1$ vectors.

By a similar argument, we now show that any affine basis of X must have size $k + 1$. One can rule out the possibility of infinite affine bases using the fact that $\dim(V) = n$ is finite. Now, let $L = (y_0, y_1, \dots, y_s)$ be any ordered affine basis of X ; we will show $s = k$ (so the list has size $k + 1$). We know $L' = (y_1 - y_0, \dots, y_s - y_0)$ is a linearly independent list. Moreover, $X = y_0 + W$, so each vector in L' lies in W . Let us check that W is the linear span of L' . Given $w \in W$, we have $y_0 + w \in X$, so that $y_0 + w = c_0 y_0 + c_1 y_1 + \cdots + c_s y_s$ for some $c_0, c_1, \dots, c_s \in F$ such that $c_0 + c_1 + \cdots + c_s = 1$. Solving for w gives

$$w = (-1 + c_0 + c_1 + \cdots + c_s)y_0 + c_1(y_1 - y_0) + \cdots + c_s(y_s - y_0).$$

The coefficient of y_0 is zero, so w is a linear combination of the vectors in L' . Thus, L' is an ordered basis of the k -dimensional space W , forcing $s = k$.

These proofs illustrate the technique of establishing facts about affine concepts involving X by looking at the corresponding linear concepts involving the direction subspace of X . As further examples of this method, we ask the reader to prove that any affinely independent subset of X can be extended to an affine basis of X , whereas any affine spanning set for X contains a subset that is an affine basis of X (Exercise 12). Furthermore, the maximum size of an affinely independent set in an n -dimensional space V is $n + 1$.

Now let X be an affine set with ordered affine basis $L = (v_0, v_1, \dots, v_k)$. We claim that for each $z \in X$, there exist unique scalars $c_0, c_1, \dots, c_k \in F$ with $\sum_{i=0}^k c_i = 1$ and $\sum_{i=0}^k c_i v_i = z$. These scalars are called the *barycentric coordinates of z relative to L* . Existence of the c_i follows since $X = \text{aff}(\{v_0, \dots, v_k\})$. To see that the c_i 's are unique, suppose we also had $d_0, d_1, \dots, d_k \in F$ with $\sum_{i=0}^k d_i = 1$ and $\sum_{i=0}^k d_i v_i = z$. Subtracting the two expressions for z gives $\sum_{i=0}^k (c_i - d_i) v_i = 0$, where the sum of the coefficients $c_i - d_i$ is $1 - 1 = 0$. By affine independence of L , $c_i - d_i = 0$ for all i , so $c_i = d_i$ for all i .

Let us work out an example of barycentric coordinates for an affine plane in \mathbb{R}^3 . Let P be the set of points $(x, y, z) \in \mathbb{R}^3$ satisfying $x + 2y - 3z = 5$. By choosing values for y and z and calculating x , we obtain the three vectors $u = (5, 0, 0)$, $v = (3, 1, 0)$, and $w = (2, 0, -1)$ in P . One verifies that (u, v, w) is an ordered affine basis for P . Consider a point $(-8, 2, -3)$ in P . We can express this point as the affine combination

$$(-8, 2, -3) = -4(5, 0, 0) + 2(3, 1, 0) + 3(2, 0, -1) = -4u + 2v + 3w.$$

So, the point in P with Cartesian coordinates $(-8, 2, -3)$ has barycentric coordinates $(-4, 2, 3)$ relative to the ordered affine basis (u, v, w) .

11.9 Characterizations of Affine Sets

We defined affine sets to be those subsets of the n -dimensional vector space V that are closed under affine combinations. As in the case of linear subspaces, we now give several other characterizations of which nonempty subsets of V are affine. We will show that, ignoring the empty set: (a) every affine set is a translate of a linear subspace of V ; (b) every affine set is the affine span of an ordered list of k affinely independent vectors, for some k between 1 and $n + 1$; (c) every $(n - m)$ -dimensional affine set in F^n is the solution set of a system of m (possibly non-homogeneous) nonzero linear equations $Ax = b$ for some $A \in M_{m,n}(F)$ and $b \in F^m$; (d) every $(n - m)$ -dimensional affine set in V is an intersection of m affine hyperplanes in V .

We proved (a) in §11.5 and (b) in §11.8. For (c), let X be an $(n - m)$ -dimensional affine set in F^n . Using (a), write $X = u + W$ for some $(n - m)$ -dimensional linear subspace of F^n and some $u \in F^n$. By (b) of §11.3, there is a matrix $A \in M_{m,n}(F)$ such that $W = \{w \in F^n : Aw = 0\}$. We claim that for $b = Au$, we have $X = \{x \in F^n : Ax = b\}$. On one hand, if $x \in X$, then $x = u + w$ for some $w \in W$, so $Ax = A(u+w) = Au + Aw = b + 0 = b$. On the other hand, if $x \in F^n$ satisfies $Ax = b$, then $A(x - u) = Ax - Au = b - b = 0$, so $x - u \in W$ and $x \in u + W = X$. We prove (d) similarly: write $X = u + W$ as in (c), and invoke (c) of §11.3 to find linear hyperplanes H_1, \dots, H_m in V with $W = H_1 \cap \dots \cap H_m$. Each translate $u + H_i$ is an affine hyperplane in V , and one checks that

$$X = u + W = u + (H_1 \cap \dots \cap H_m) = (u + H_1) \cap \dots \cap (u + H_m).$$

11.10 Affine Maps

Let V and W be two vector spaces over the field F . A *linear map* or *linear transformation* is a function $T : V \rightarrow W$ such that $T(x + y) = T(x) + T(y)$ and $T(cx) = cT(x)$ for all $x, y \in V$ and all $c \in F$. From this definition and induction, one sees that a linear map T satisfies $T(c_1x_1 + \cdots + c_kx_k) = c_1T(x_1) + \cdots + c_kT(x_k)$ for all $k \in \mathbb{N}$, $c_i \in F$, and $x_i \in V$. In other words, *linear maps preserve linear combinations*. By analogy, we define a map $U : V \rightarrow W$ to be an *affine map* or *affine transformation* iff U preserves affine combinations, i.e., $U(c_1x_1 + \cdots + c_kx_k) = c_1U(x_1) + \cdots + c_kU(x_k)$ for all $k \in \mathbb{N}$, $x_i \in V$, and $c_i \in F$ satisfying $c_1 + \cdots + c_k = 1$. (When $F = \mathbb{R}$, it suffices to check the condition for $k = 2$; see Exercise 26.) More generally, we allow the domain (or codomain) of an affine map to be an affine subset of V (or W). An *affine isomorphism* is a bijective affine map.

Certainly, every linear map is an affine map, but the converse is not true. For example, given a fixed $u \in V$, the *translation map* $T_u : V \rightarrow V$ defined by $T_u(v) = u + v$ for $v \in V$ is an affine map. To check this, fix $k \in \mathbb{N}$, $x_i \in V$, and $c_i \in F$ with $\sum_{i=1}^k c_i = 1$, and compute

$$\begin{aligned} T_u(c_1x_1 + \cdots + c_kx_k) &= u + c_1x_1 + \cdots + c_kx_k = (c_1 + \cdots + c_k)u + c_1x_1 + \cdots + c_kx_k \\ &= c_1(u + x_1) + \cdots + c_k(u + x_k) = c_1T_u(x_1) + \cdots + c_kT_u(x_k). \end{aligned}$$

However, since linear maps must send zero to zero, T_u is a linear map only when $u = 0$. More generally, the reader can check that for $A \in M_{m,n}(F)$ and $b \in F^m$, the map $U : F^n \rightarrow F^m$ given by $U(x) = Ax + b$ for $x \in F^n$ is always affine, but is linear only when $b = 0$.

In fact, a map $U : V \rightarrow W$ is a linear map iff U is an affine map and $U(0_V) = 0_W$. The forward implication is immediate from the definitions. Conversely, assume U is an affine map with $U(0) = 0$. Given $x, y \in V$ and $c \in F$, compute:

$$\begin{aligned} U(x + y) &= U(-1 \cdot 0 + 1x + 1y) = -1U(0) + 1U(x) + 1U(y) = U(x) + U(y); \\ U(cx) &= U((1 - c) \cdot 0 + cx) = (1 - c)U(0) + cU(x) = cU(x). \end{aligned}$$

So U is linear.

The following facts are routinely verified: the identity map on any affine set is an affine isomorphism; the composition of two affine maps (resp. affine isomorphisms) is an affine map (resp. affine isomorphism); and the inverse of an affine isomorphism is also an affine isomorphism. For example, every translation map T_u is an affine isomorphism with inverse $T_u^{-1} = T_{-u}$, which is also affine. One can also check that the direct or inverse image of an affine set under an affine map is an affine set.

We know that every linear map $T : F^n \rightarrow F^m$ has the form $T(x) = Ax$ ($x \in F^n$) for a unique matrix $A \in M_{m,n}(F)$ (namely, A is the matrix whose j 'th column is $T(e_j)$ for $1 \leq j \leq n$). Using this, we show that every affine map $U : F^n \rightarrow F^m$ has the form $U(x) = Ax + b$ for a unique $A \in M_{m,n}(F)$ and $b \in F^m$. Fix the affine map U . Choose $b = U(0) \in F^m$, and consider the map $T_b^{-1} \circ U : F^n \rightarrow F^m$, where T_b^{-1} is the inverse of translation by b in the space F^m . The composition $T_b^{-1} \circ U$ is an affine map sending zero to zero, so it is a linear map from F^n to F^m . Thus, there is a matrix A with $T_b^{-1} \circ U(x) = Ax$ for all $x \in F^n$. Solving for U , we see that $U(x) = T_b(Ax) = Ax + b$ for all $x \in F^n$. Now, b is uniquely determined by U since $U(x) = Ax + b$ forces $b = A0 + b = U(0)$ no matter what A is. Knowing this, A is uniquely determined because it is the unique matrix representing the linear map $T_b^{-1} \circ U$. This completes the proof. More generally, for abstract vector spaces V and W , a similar argument shows that every affine map $U : V \rightarrow W$ has the form $U = T_b \circ S$ for a uniquely determined linear map $S : V \rightarrow W$ and a uniquely determined $b \in W$ (here, T_b is translation by b on W).

The following *universal mapping property (UMP)* for affine bases lets us build affine maps by specifying how the map should act on an affine basis. Suppose $X \subseteq V$ is an affine set with affine basis $S = \{v_0, \dots, v_k\}$. For any $y_0, \dots, y_k \in W$, there exists a unique affine map $U : X \rightarrow W$ with $U(v_i) = y_i$ for $0 \leq i \leq k$. To give a specific formula for U , note that each $x \in X$ has a unique expression as an affine combination $x = \sum_{i=0}^k c_i v_i$ with $c_i \in F$ and $\sum_{i=0}^k c_i = 1$. So if U exists at all, it must be given by the formula $U(x) = \sum_{i=0}^k c_i U(v_i) = \sum_{i=0}^k c_i y_i$. A routine calculation confirms that this map is affine and sends each v_i to y_i . One can also prove the UMP by reducing to the direction subspace of X and appealing to the universal mapping property for the linear basis $\{v_1 - v_0, \dots, v_k - v_0\}$ of this subspace (Exercise 25).

11.11 Convex Sets

For our study of convexity in the rest of this chapter, we shall consider only real vector spaces, taking $V = \mathbb{R}^n$. To motivate the definition of a convex set, consider once again the line L shown in Figure 11.1. The entire line L through v and w is obtained by taking affine combinations $(1-t)v + tw$ as t ranges over the entire real line \mathbb{R} . In contrast, if we only wanted to get the points of L on the line segment joining v to w , we would only take the affine combinations $(1-t)v + tw = v + t(w-v)$ with t ranging through real numbers in the closed interval $I = [0, 1]$. (Taking $t = 0$ gives v , taking $t = 1$ gives w , and intermediate t 's give precisely the points in the interior of the line segment.) In general, for any $v, w \in \mathbb{R}^n$, we define the *line segment joining v and w* to be the set $\{(1-t)v + tw : t \in \mathbb{R}, 0 \leq t \leq 1\}$.

A subset C of \mathbb{R}^n is *convex* iff for all $v, w \in C$ and all $t \in \mathbb{R}$ with $0 \leq t \leq 1$, $(1-t)v + tw$ lies in C . Geometrically, a set C is convex iff the line segment joining any two points of C is always contained in C . All affine sets are convex, but the converse is not true. For example, any finite closed interval $[a, b]$ in \mathbb{R}^1 is readily seen to be convex but not affine.

Given finitely many points $v_1, \dots, v_k \in \mathbb{R}^n$, a *convex combination* of these points is a linear combination $c_1 v_1 + \dots + c_k v_k$ with $c_1 + \dots + c_k = 1$ and $c_i \geq 0$ for all i (and hence $c_i \in [0, 1]$ for all i). Let us show that a convex set $C \subseteq \mathbb{R}^n$ is closed under convex combinations; i.e., for all $k \in \mathbb{N}^+$ and all $v_1, \dots, v_k \in C$, every convex combination of the v_i 's lies in C . We argue by induction on k . If $k = 1$, we must have $c_1 = 1$. Given that $v_1 \in C$, we certainly have $c_1 v_1 = 1v_1 = v_1 \in C$. If $k = 2$, we must have $c_1 = 1 - c_2$, so $c_1 v_1 + c_2 v_2 \in C$ by definition of a convex set. Now, fix $k \geq 3$, and assume that convex combinations of $k - 1$ or fewer points of C are already known to lie in C . Fix $v_1, \dots, v_k \in C$ and $c_1, \dots, c_k \in [0, 1]$ with $\sum_{i=1}^k c_i = 1$. If $c_1 = 1$, then the conditions on the c_i 's force $c_2 = \dots = c_k = 0$, so that $c_1 v_1 + \dots + c_k v_k = v_1 \in C$. Otherwise, let $d_i = c_i / (1 - c_1)$ for $2 \leq i \leq k$, and note that

$$c_1 v_1 + c_2 v_2 + \dots + c_k v_k = c_1 v_1 + (1 - c_1)(d_2 v_2 + \dots + d_k v_k).$$

Since c_2, \dots, c_k are nonnegative real numbers with sum $1 - c_1$, it follows that d_2, \dots, d_k are nonnegative real numbers with sum 1. Then $w = d_2 v_2 + \dots + d_k v_k$ is a convex combination of $k - 1$ elements of C , which is in C by the induction hypothesis. So $\sum_{i=1}^k c_i v_i = (1 - c_1)w + c_1 v_1 \in C$ by definition of a convex set. This completes the induction proof.

11.12 Convex Hulls

Given $S \subseteq \mathbb{R}^n$, we have studied the *linear span* of S , which consists of all linear combinations of elements of S , and the *affine span* of S , which consists of all affine combinations of elements of S . By analogy, we define the *convex span* or *convex hull* of $S \subseteq \mathbb{R}^n$, denoted $\text{conv}(S)$, to be the set of all convex combinations of finitely many elements of S .

We assert that $\text{conv}(S)$ is a convex set containing S , and $\text{conv}(S) \subseteq T$ for all convex sets T containing S . (So, we could have equally well defined the convex hull of S to be the intersection of all convex subsets of \mathbb{R}^n containing S .) To check convexity of $\text{conv}(S)$, fix $x, y \in \text{conv}(S)$ and $t \in [0, 1]$; we must check that $z = tx + (1-t)y \in \text{conv}(S)$. By definition, there are $v_1, \dots, v_k \in S$, $c_1, \dots, c_k \in [0, 1]$ with sum 1, and $d_1, \dots, d_k \in [0, 1]$ with sum 1, such that $x = \sum_{i=1}^k c_i v_i$ and $y = \sum_{i=1}^k d_i v_i$. (By taking some c_i 's and d_i 's equal to zero, we can assume that the same elements v_1, \dots, v_k of S are used in the expressions for x and y .) Now compute

$$z = tx + (1-t)y = \sum_{i=1}^k (tc_i + (1-t)d_i)v_i.$$

Each c_i and d_i is a nonnegative real number, as are t and $1-t$, so each coefficient $tc_i + (1-t)d_i$ is a nonnegative real number. The sum of these coefficients is

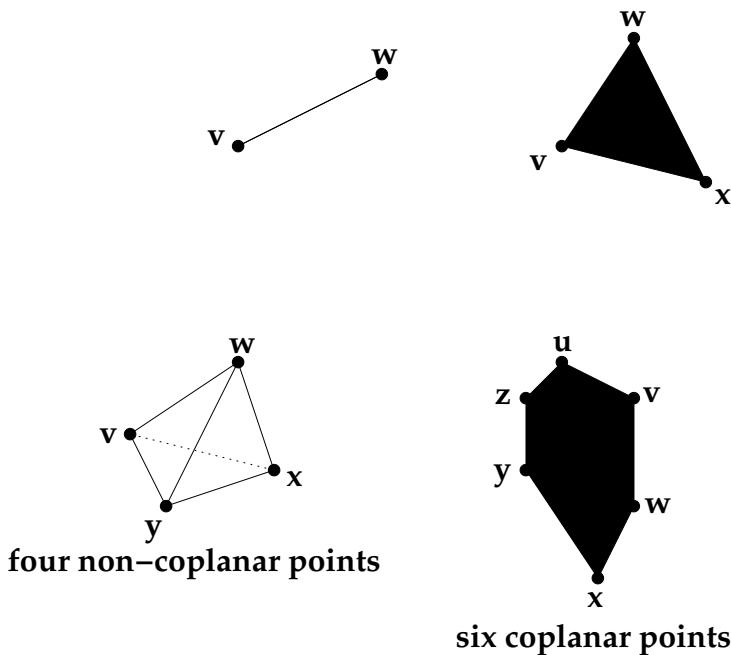
$$\sum_{i=1}^k (tc_i + (1-t)d_i) = t \sum_{i=1}^k c_i + (1-t) \sum_{i=1}^k d_i = t \cdot 1 + (1-t) \cdot 1 = 1.$$

So z is a convex combination of $v_1, \dots, v_k \in S$, hence $z \in \text{conv}(S)$ as required. Since any $v \in S$ is the convex combination $v = 1v$, we have $S \subseteq \text{conv}(S)$. Finally, if T is any convex subset of \mathbb{R}^n containing S , we know $\text{conv}(S) \subseteq T$ since T is closed under convex combinations of its elements.

The convex hulls of some finite sets of points are illustrated in Figure 11.2. In general, the convex hull of a one-point set $\{v\}$ is $\{v\}$; the convex hull of $\{v, w\}$ is the line segment joining v and w ; the convex hull of three non-collinear points $\{v, w, x\}$ is the triangle (including interior) with vertices v, w, x ; and the convex hull of four non-coplanar points $\{v, w, x, y\}$ is the tetrahedron (including interior) with vertices v, w, x, y . For $k \geq 0$, we define a k -dimensional simplex to be a set $\text{conv}(\{v_0, v_1, \dots, v_k\})$, where $\{v_0, v_1, \dots, v_k\}$ is an affinely independent subset of \mathbb{R}^n . Geometrically, simplices are k -dimensional analogues of triangles and tetrahedra. A specific example of a k -dimensional simplex is $\Delta_k = \text{conv}(\{0, e_1, \dots, e_k\})$, where the e_i 's are standard basis vectors in \mathbb{R}^k . Writing out the definition, we see that Δ_k is the set of vectors $(c_1, \dots, c_k) \in \mathbb{R}^k$ with $c_i \geq 0$ for all i and $\sum_{i=1}^k c_i \leq 1$. (This sum can be less than 1, since the scalar c_0 disappears in the convex combination $c_0 \cdot 0 + c_1 e_1 + \dots + c_k e_k$.)

11.13 Carathéodory's Theorem

Consider the six-point set $S = \{u, v, w, x, y, z\} \subseteq \mathbb{R}^2$ shown in the lower-right part of Figure 11.2. The convex hull $\text{conv}(S)$ consists of all convex combinations of the six points of S . If we tried to create this convex hull by drawing line segments between every pair of points in S , the resulting set would not be convex (it would consist of the sides and diagonals of the shaded hexagonal region). On the other hand, suppose we took all convex

**FIGURE 11.2**

Examples of Convex Hulls.

combinations of three-element subsets of S . For each fixed subset of size 3, the set of convex combinations using these three elements is a triangle with interior. One sees from the figure that the union of these triangles is all of $\text{conv}(S)$. Thus, to obtain the convex hull in this example, it suffices to use only convex combinations of three or fewer points from S .

Carathéodory's theorem generalizes this remark to convex hulls in higher dimensions. The theorem asserts that *for any subset S of \mathbb{R}^n , $\text{conv}(S)$ consists of all convex combinations of $n + 1$ or fewer points of S .* To prove this, fix $S \subseteq \mathbb{R}^n$ and pick any point $w \in \text{conv}(S)$. By definition, we know there exist $m \in \mathbb{N}$, points $v_0, v_1, \dots, v_m \in S$, and nonnegative real scalars c_0, c_1, \dots, c_m such that $c_0 + c_1 + \dots + c_m = 1$ and $w = c_0v_0 + c_1v_1 + \dots + c_mv_m$. If $m \leq n$, the claimed result holds for the point w . So assume $m > n$. We will show how to find another expression for w as a convex combination of fewer than $m + 1$ points of S . Continuing this reduction process for at most $m - n$ steps, we eventually arrive at an expression for w as a convex combination of at most $n + 1$ points of S .

We can assume all v_i 's are distinct, since otherwise we can reduce m by using $cv + dv = (c + d)v$ to group together two terms in the convex combination. We can assume all $c_i > 0$, since otherwise we can reduce m by deleting the term $c_i v_i = 0$ from the sum. Now, since $m + 1 > n + 1$, the set $\{v_0, \dots, v_m\}$ in \mathbb{R}^n must be affinely dependent. So there exist $d_0, \dots, d_m \in \mathbb{R}$ with some $d_i \neq 0$, $d_0 + \dots + d_m = 0$, and $d_0v_0 + \dots + d_mv_m = 0$. At least one d_i is strictly positive, and we can reorder terms to ensure that $d_m > 0$ and $c_m/d_m \leq c_i/d_i$ for all $i < m$ such that $d_i > 0$. Now define $b_i = c_i - (c_m/d_m)d_i$ for $0 \leq i < m$. We claim each $b_i \geq 0$. If $d_i > 0$, the claim follows by multiplying $c_m/d_m \leq c_i/d_i$ by the positive quantity d_i and rearranging. If $d_i \leq 0$, the claim follows since $c_i > 0$, $c_m/d_m > 0$,

hence $-(c_m/d_m)d_i \geq 0$ and $b_i \geq 0$. We also claim the sum of the b_i 's is 1. For,

$$\sum_{i=0}^{m-1} b_i = \sum_{i=0}^{m-1} c_i - (c_m/d_m) \sum_{i=0}^{m-1} d_i = (1 - c_m) - (c_m/d_m)(-d_m) = 1.$$

Finally, we claim that $w = \sum_{i=0}^{m-1} b_i v_i$. For,

$$\sum_{i=0}^{m-1} b_i v_i = \sum_{i=0}^{m-1} c_i v_i - (c_m/d_m) \sum_{i=0}^{m-1} d_i v_i = (w - c_m v_m) - (c_m/d_m)(-d_m v_m) = w.$$

We have now expressed w as a convex combination of m points of S , which completes the proof.

Carathéodory's theorem can be used to prove that the convex hull of a closed and bounded (compact) subset of \mathbb{R}^n is closed and bounded (Exercise 45).

11.14 Hyperplanes and Half-Spaces in \mathbb{R}^n

For every affine hyperplane H in \mathbb{R}^n , there exist $a_1, \dots, a_n, b \in \mathbb{R}$ such that

$$H = \{(x_1, \dots, x_n) \in \mathbb{R}^n : a_1 x_1 + \dots + a_n x_n = b\};$$

this is a special case of part (c) of the theorem in §11.9. Recall that the *dot product* or *inner product* of two vectors $w = (w_1, \dots, w_n)$ and $x = (x_1, \dots, x_n)$ in \mathbb{R}^n is defined by

$$w \bullet x = \langle w, x \rangle = w_1 x_1 + w_2 x_2 + \dots + w_n x_n.$$

Letting $a = (a_1, \dots, a_n) \neq 0$, we can write $H = \{x \in \mathbb{R}^n : a \bullet x = b\}$. Recall that the *Euclidean length* of the vector a is $\|a\| = \sqrt{a \bullet a} = (a_1^2 + a_2^2 + \dots + a_n^2)^{1/2}$. By multiplying both sides of the equation $a \bullet x = b$ by the nonzero scalar $1/\|a\|$, we obtain an equivalent equation $u \bullet x = c$ where u is a *unit vector* (i.e., $\|u\| = 1$) and $c = \|a\|^{-1}b \in \mathbb{R}$. By multiplying by -1 if needed, we can arrange that $c \geq 0$. When $c \neq 0$, one can check that the unit vector u and positive scalar c are uniquely determined by H . If $c = 0$, there are exactly two possible unit vectors u that can be used in the equation defining H .

Given an affine hyperplane H , write $H = \{x \in \mathbb{R}^n : u \bullet x = c\}$ with $\|u\| = 1$ and $c \geq 0$. If $c = 0$, H consists of all vectors $x \in \mathbb{R}^n$ perpendicular to the unit vector u . If $c > 0$ and x_0 is any fixed point on H , $x \in H$ iff $u \bullet x = u \bullet x_0$ iff $u \bullet (x - x_0) = 0$ iff $x - x_0$ is perpendicular to the unit vector u . Accordingly, u and its nonzero scalar multiples are called *normal vectors* for H .

The two *closed half-spaces* determined by H are the sets $\{x \in \mathbb{R}^n : u \bullet x \geq c\}$ and $\{x \in \mathbb{R}^n : u \bullet x \leq c\}$. Geometrically, the first set consists of points of H together with points outside of H “on the same side of H as u ” (where u is a vector with tail drawn at some point $x_0 \in H$), and the second set consists of H together with the points on the opposite side of H from u . Similarly, the two *open half-spaces* determined by H are the sets $\{x \in \mathbb{R}^n : u \bullet x > c\}$ and $\{x \in \mathbb{R}^n : u \bullet x < c\}$. All four half-spaces determined by H , along with H itself, are convex. For example, suppose $x, y \in S = \{z \in \mathbb{R}^n : u \bullet z < c\}$ and $t \in [0, 1]$. We know $u \bullet x < c$ and $u \bullet y < c$. Since the inner product on \mathbb{R}^n is linear in each variable and $0 \leq t \leq 1$, we have

$$u \bullet (tx + (1-t)y) = t(u \bullet x) + (1-t)(u \bullet y) < tc + (1-t)c = c,$$

so $tx + (1 - t)y \in S$. A closed half-space is the solution set of a *linear inequality* $u_1x_1 + \dots + u_nx_n \leq c$, and similarly an open half-space is the solution set of a *strict linear inequality*. It is a routine exercise to check that the intersection of any family of convex subsets of \mathbb{R}^n is also convex. Taking each convex set in the family to be a hyperplane or a closed half-space or an open half-space, we see that *the solution set of any family (possibly infinite) of linear equations, inequalities, and strict inequalities in n variables is a convex subset of \mathbb{R}^n* . Observe that each linear equation $u \bullet x = c$ appearing in the family could be replaced by the two linear inequalities $u \bullet x \leq c$ and $(-u) \bullet x \leq c$.

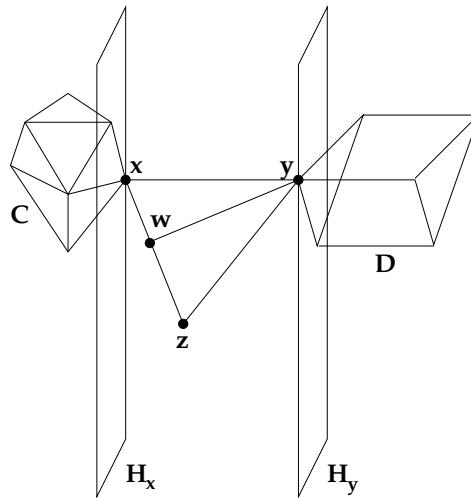
11.15 Closed Convex Sets

We would like to be able to describe convex sets as intersections of closed half-spaces, just as we described affine sets as intersections of affine hyperplanes. However, not all convex sets can be written in this form. To see why, recall that a subset C of \mathbb{R}^n is *closed* iff for every sequence of vectors v_1, \dots, v_k, \dots in C that converge to a vector $v \in \mathbb{R}^n$ (i.e., $\lim_{k \rightarrow \infty} \|v_k - v\| = 0$), the limit vector v lies in C . In words, C is closed (in the topological sense) iff C is closed under taking limits of convergent sequences of vectors in C . Using this definition, one can check quickly that the intersection of any family of closed subsets of \mathbb{R}^n is also a closed set. Furthermore, every closed half-space is a closed set (as the name suggests). Combining these facts, we see that any intersection of closed half-spaces will always be a closed (and convex) set. We want to prove that, conversely, *every closed and convex subset of \mathbb{R}^n is the intersection of all closed half-spaces containing it*.

To obtain this result, we will use a theorem that lets us “separate” certain subsets of \mathbb{R}^n using an affine hyperplane. Recall that a set $D \subseteq \mathbb{R}^n$ is *bounded* iff there exists $M \in \mathbb{R}$ such that for all $x \in D$, $\|x\| \leq M$. Our *separation theorem* says that *for all disjoint, closed, convex sets C and D in \mathbb{R}^n with C bounded, there is an affine hyperplane H in \mathbb{R}^n (with associated open half-spaces denoted H^+ and H^-) such that $C \subseteq H^+$ and $D \subseteq H^-$ or vice versa.*

To start the proof, fix sets C and D satisfying the hypotheses of the theorem. For points $x, y \in \mathbb{R}^n$, let $d(x, y) = \|x - y\| = \sqrt{(x - y) \bullet (x - y)}$ be the *distance* between x and y in \mathbb{R}^n . For a point $x \in \mathbb{R}^n$ and a nonempty set $B \subseteq \mathbb{R}^n$, let $d(x, B) = \inf_{y \in B} d(x, y)$ be the *distance between x and B* . For two nonempty sets A and B in \mathbb{R}^n , let $d(A, B) = \inf_{x \in A} d(x, B) = \inf_{x \in A, y \in B} d(x, y)$ be the *distance between A and B* . Now, it is shown in real analysis that a closed bounded subset C of \mathbb{R}^n is *sequentially compact*; i.e., every sequence of points of C has a convergent subsequence. Using compactness, one can prove that every continuous function $f : C \rightarrow \mathbb{R}$ with compact domain *attains its minimum value*; i.e., there is $x \in C$ with $f(x) \leq f(w)$ for all $w \in C$. With these definitions and facts in hand, one may check that there exist $x \in C$ and $y \in D$ with $d(x, y) = d(C, D)$ (see Exercise 40). Set $r = d(x, y) = \|y - x\|$; since C and D are disjoint, $x \neq y$ and so $r > 0$.

Let u be the unit vector $r^{-1}(y - x)$ parallel to the line segment from x to y in \mathbb{R}^n . Let $c_1 = u \bullet x$ and $c_2 = u \bullet y$. Finally, let $H_x = \{v \in \mathbb{R}^n : u \bullet v = c_1\}$ be the affine hyperplane through x perpendicular to u , and let $H_y = \{v \in \mathbb{R}^n : u \bullet v = c_2\}$ be the affine hyperplane through y perpendicular to u . We cannot have $c_1 = c_2$, since otherwise $u \bullet (y - x) = c_2 - c_1 = 0$ would give $r = r(u \bullet u) = u \bullet (ru) = 0$. We consider the case $c_1 < c_2$ below; the case $c_1 > c_2$ is handled similarly. We claim that $C \subseteq \{v \in \mathbb{R}^n : u \bullet v \leq c_1\}$ and that $D \subseteq \{v \in \mathbb{R}^n : u \bullet v \geq c_2\}$. The claim is illustrated by Figure 11.3, in which C is contained in the closed half-space weakly left of H_x , and D is contained in the closed

**FIGURE 11.3**Finding a Hyperplane to Separate C and D .

half-space weakly right of H_y . The separation theorem will follow from the claim by taking $H = \{z \in \mathbb{R}^n : u \bullet z = c_3\}$ for any choice of $c_3 \in \mathbb{R}$ with $c_1 < c_3 < c_2$.

We now prove the claim. Fix a point $z \in \mathbb{R}^n$ with $u \bullet z = c > c_1 = u \bullet x$; we will show $z \notin C$. Suppose $z \in C$ to obtain a contradiction. Since $x, z \in C$ and C is convex, the entire line segment joining x and z is contained in C . Figure 11.3 suggests that we can drop an altitude from y to this line segment to find a point $w \in C$ with $d(w, y) < d(x, y) = d(C, D)$, giving a contradiction. To make this pictorial argument precise, we consider the function $g : [0, 1] \rightarrow \mathbb{R}$ such that $g(t) = d(y, x + t(z - x))^2$ for $t \in [0, 1]$. By definition of distance and the linearity of the dot product in each argument, we know

$$\begin{aligned} g(t) &= [(y - x) - t(z - x)] \bullet [(y - x) - t(z - x)] \\ &= (y - x) \bullet (y - x) - 2t(y - x) \bullet (z - x) + t^2(z - x) \bullet (z - x). \end{aligned}$$

Recalling $r = d(x, y)$, $y - x = ru$, and setting $s = d(z, x)$, this becomes

$$g(t) = r^2 - 2t((ru) \bullet z - (ru) \bullet x) + s^2t^2 = s^2t^2 - 2r(c - c_1)t + r^2.$$

The graph of g in \mathbb{R}^2 is a parabola with slope $g'(0) = -2r(c - c_1) < 0$ at $t = 0$. So, for a sufficiently small positive $t \in [0, 1]$, $w = x + t(z - x)$ will satisfy $d(w, y) = \sqrt{g(t)} < \sqrt{g(0)} = d(y, x)$. This contradiction shows $z \notin C$. A similar argument proves that no point $z \in \mathbb{R}^n$ with $u \bullet z < c_2$ can lie in D . This finishes the proof of the claim and the separation theorem.

We can now prove that a closed convex set D in \mathbb{R}^n is the intersection E of all closed half-spaces S with $D \subseteq S$. The inclusion $D \subseteq E$ is evident from the definition of E . For the reverse inclusion, suppose $x \in \mathbb{R}^n$ and $x \notin D$. Apply the separation theorem to the one-point set $C = \{x\}$ and the closed convex set D ; note that C is closed, convex, bounded, and disjoint from D . We obtain a hyperplane H with $C \subseteq H^+$ (one of the open half-spaces for H) and $D \subseteq H^-$ (the other open half-space for H). Then $S = H^- \cup H$ is a closed half-space that contains D but not x . This is one of the half-spaces we intersect to obtain E , so $x \notin E$ as needed. In fact, the proof shows that D is also the intersection of all the open half-spaces containing D .

11.16 Cones and Convex Cones

We would like to understand the structure of subsets of \mathbb{R}^n formed by intersecting a *finite* number of closed half-spaces. Our ultimate goal is to prove that *a subset S of \mathbb{R}^n is the convex hull of a finite set of points iff S is bounded and $S = H_1 \cap \dots \cap H_k$ for finitely many closed half-spaces H_i* . This theorem is geometrically very plausible for subsets of \mathbb{R}^2 and \mathbb{R}^3 (cf. Figures 11.2 and 11.3), but it is quite tricky to prove in general dimensions. To prove this result, and to describe unbounded intersections of finitely many closed half-spaces, we must first develop some machinery involving cones.

A subset C of \mathbb{R}^n is a *cone* iff $0 \in C$ and for all $x \in C$ and all $b > 0$ in \mathbb{R} , $bx \in C$. In other words, cones are sets containing 0 that are closed under multiplication by *positive* real scalars. Equivalently, cones are nonempty sets closed under multiplication by nonnegative scalars. We claim *a cone C is convex iff C is closed under addition*. On one hand, suppose C is a convex cone and $x, y \in C$. Then $z = (1/2)x + (1/2)y \in C$ by convexity, so $x + y = 2z \in C$ because C is a cone, so C is closed under addition. On the other hand, suppose a cone C is closed under addition. For $x, y \in C$ and $t \in [0, 1]$, we know $tx \in C$ and $(1-t)y \in C$, since $t \geq 0$ and $1-t \geq 0$. Then $tx + (1-t)y \in C$, so that C is convex.

A *positive combination* (or *conical combination*) of vectors $v_1, \dots, v_k \in \mathbb{R}^n$ is a linear combination $c_1v_1 + \dots + c_kv_k$, where all $c_i \in \mathbb{R}$ are positive or zero. By the remarks in the previous paragraph and induction on k , we see that *convex cones are closed under positive combinations of their elements*. Conversely, any subset of \mathbb{R}^n closed under positive combinations is evidently a convex cone (taking $k = 0$ shows that such a subset must contain zero).

Given $S \subseteq \mathbb{R}^n$, the *convex cone generated by S* is the set $\text{cone}(S)$ consisting of all positive combinations of vectors $v_1, \dots, v_k \in S$. Let us show that $\text{cone}(S)$ is a convex cone containing S . Taking $k = 0$, we see that $0 \in \text{cone}(S)$. Fix $x, y \in \text{cone}(S)$ and $b > 0$ in \mathbb{R} . We can write $x = c_1v_1 + \dots + c_kv_k$ and $y = d_1v_1 + \dots + d_kv_k$ for some real scalars $c_i, d_i \geq 0$ and some $v_1, \dots, v_k \in S$. Then $bx = \sum_{i=1}^k(bc_i)v_i \in \text{cone}(S)$ since each $bc_i \geq 0$, and $x + y = \sum_{i=1}^k(c_i + d_i)v_i \in \text{cone}(S)$ since each $c_i + d_i \geq 0$. So $\text{cone}(S)$ is a convex cone. This cone contains S since each $v \in S$ is the positive combination $v = 1v$. Moreover, if T is any convex cone containing S , then $\text{cone}(S) \subseteq T$ since T is closed under positive combinations of its elements. We conclude that $\text{cone}(S)$ can also be described as the intersection of all convex cones in \mathbb{R}^n containing S .

All linear subspaces W of \mathbb{R}^n are convex cones, as they are closed under all linear combinations of their elements (hence are closed under positive combinations). Define a *linear half-space* of \mathbb{R}^n to be any closed half-space determined by a linear hyperplane in \mathbb{R}^n . Explicitly, linear half-spaces are sets of the form $\{x \in \mathbb{R}^n : u \bullet x \leq 0\}$ for some nonzero $u \in \mathbb{R}^n$. One checks immediately that linear half-spaces are convex cones. Moreover, any intersection of convex cones is also a convex cone. In particular, intersections of linear half-spaces are convex cones.

By imposing finiteness conditions on the constructions in the two preceding paragraphs, we obtain two special classes of convex cones. First, a *V-cone* is a convex cone of the form $\text{cone}(S)$, where S is a *finite* subset of \mathbb{R}^n . Explicitly, $C \subseteq \mathbb{R}^n$ is a V-cone iff there exist $k \in \mathbb{N}^+$ and $v_1, \dots, v_k \in \mathbb{R}^n$ such that

$$C = \{c_1v_1 + \dots + c_kv_k : c_i \in \mathbb{R}, c_i \geq 0 \text{ for } 1 \leq i \leq k\}.$$

We call v_1, \dots, v_k *generators* of the cone C . Second, an *H-cone* is an intersection of finitely many linear half-spaces in \mathbb{R}^n . Explicitly, $D \subseteq \mathbb{R}^n$ is an H-cone iff there exist $k \in \mathbb{N}^+$ and

$u_1, \dots, u_k \in \mathbb{R}^n$ such that

$$D = \{x \in \mathbb{R}^n : u_i \bullet x \leq 0 \text{ for } 1 \leq i \leq k\}.$$

The next four sections are devoted to proving that *a subset of \mathbb{R}^n is an H-cone iff it is a V-cone*. This theorem will be the key to understanding the structure of finite intersections of general closed half-spaces.

11.17 Intersection Lemma for V-Cones

One can readily check that the intersection of any H-cone in \mathbb{R}^n with any linear subspace W is also an H-cone. To help prove that H-cones are the same as V-cones, we first need to prove that V-cones satisfy a special case of this *intersection property*. Specifically, given a V-cone C in \mathbb{R}^n , let W be the linear subspace $\mathbb{R}^{n-1} \times \{0\}$, which is the solution set of the equation $x_n = 0$. We will prove that $C \cap W$ is also a V-cone.

Since C is a V-cone, we have $C = \text{cone}(\{v_1, \dots, v_m\})$ for some $m \in \mathbb{N}^+$ and $v_i \in \mathbb{R}^n$. Here and below, we will write $v_i(n)$ for the n 'th coordinate of the vector v_i . We want to divide the generators v_i into classes based on the sign of their last coordinates. To do this, let $I = \{1, 2, \dots, m\}$, and define

$$I_0 = \{i \in I : v_i(n) = 0\}; \quad I_+ = \{i \in I : v_i(n) > 0\}; \quad I_- = \{i \in I : v_i(n) < 0\}.$$

Next, define a V-cone

$$D = \text{cone}(\{v_i : i \in I_0\} \cup \{v_i(n)v_j - v_j(n)v_i : i \in I_+, j \in I_-\}).$$

It will suffice to show that $C \cap W = D$.

First, note that each generator of D has n 'th coordinate zero, hence lies in W . This is true by definition for the v_i 's with $i \in I_0$. Given $i \in I_+$ and $j \in I_-$, the n 'th coordinate of $v_i(n)v_j - v_j(n)v_i$ is $v_i(n)v_j(n) - v_j(n)v_i(n) = 0$. Since all generators of D lie in the subspace W , any positive combination of those generators also lies in W . Hence, $D \subseteq W$. Next, note that each generator of D is a positive combination of generators of the convex cone C (e.g., for $i \in I_+$ and $j \in I_-$, $v_i(n) \geq 0$ and $-v_j(n) \geq 0$). Since all generators of D lie in the convex cone C , the convex cone D is contained in C . In summary, $D \subseteq C \cap W$.

For the reverse inclusion, fix $v \in C \cap W$. Since $v \in C$, we can write $v = c_1v_1 + \dots + c_mv_m$ for some real scalars $c_i \geq 0$. Since $v \in W$, $0 = v(n) = c_1v_1(n) + \dots + c_mv_m(n)$. Consider two cases. Case 1: for all $i \in I$, $c_iv_i(n) = 0$. Then for all $i \in I$, $c_i = 0$ or $v_i(n) = 0$, so that the only v_i 's appearing with nonzero coefficients in the expression for v are v_i 's with $i \in I_0$. So v is a positive combination of some of the generators of D , hence $v \in D$. Case 2: for some $i \in I$, $c_iv_i(n) \neq 0$. In the equation $\sum_{i \in I} c_iv_i(n) = 0$, drop terms indexed by $i \in I_0$, and move terms indexed by $i \in I_-$ to the right side. We obtain a strictly positive quantity

$$b = \sum_{i \in I_+} c_iv_i(n) = \sum_{j \in I_-} c_j(-v_j(n)) > 0.$$

Now, using both formulas for b , compute:

$$\begin{aligned} v &= \sum_{i \in I} c_i v_i = \sum_{i \in I_0} c_i v_i + \frac{1}{b} \sum_{i \in I_+} b c_i v_i + \frac{1}{b} \sum_{j \in I_-} b c_j v_j \\ &= \sum_{i \in I_0} c_i v_i + \frac{1}{b} \sum_{i \in I_+} \sum_{j \in I_-} c_j (-v_j(n)) c_i v_i + \frac{1}{b} \sum_{j \in I_-} \sum_{i \in I_+} c_i v_i(n) c_j v_j \\ &= \sum_{i \in I_0} c_i v_i + \sum_{i \in I_+} \sum_{j \in I_-} \frac{c_i c_j}{b} (v_i(n) v_j - v_j(n) v_i). \end{aligned}$$

Note the coefficients c_i and $c_i c_j / b$ are all nonnegative, so this expression is a positive combination of the generators for D . Therefore, $v \in D$ in case 2, completing the proof that $C \cap W \subseteq D$.

We have now proved that for a V-cone C in \mathbb{R}^n , the intersection $C \cap (\mathbb{R}^{n-1} \times \{0\})$ is a V-cone in \mathbb{R}^{n-1} . Iterating this result, we deduce that *for a V-cone C in \mathbb{R}^n and any $k \leq n$, $C \cap (\mathbb{R}^{n-k} \times \{0\}^k)$ is a V-cone in \mathbb{R}^{n-k}* (where we identify \mathbb{R}^{n-k} with $\mathbb{R}^{n-k} \times \{0\}^k \subseteq \mathbb{R}^n$).

11.18 All H-Cones Are V-Cones

We now use the intersection lemma for V-cones to prove that *every H-cone is a V-cone*. Given an H-cone $C \subseteq \mathbb{R}^n$, choose $k \in \mathbb{N}^+$ and $u_1, \dots, u_k \in \mathbb{R}^n$ with

$$C = \{x \in \mathbb{R}^n : u_i \bullet x \leq 0 \text{ for } 1 \leq i \leq k\}.$$

The first step is to convert this H-cone into a V-cone in the higher-dimensional space \mathbb{R}^{n+k} . To do this, define

$$D = \{(x, y) \in \mathbb{R}^n \times \mathbb{R}^k : u_i \bullet x \leq y(i) \text{ for } 1 \leq i \leq k\}. \quad (11.2)$$

(As above, $y(i)$ denotes the i 'th component of the vector y .) It is readily checked that $0 \in D$, and D is closed under addition and positive scalar multiples, so D is a convex cone. For $1 \leq j \leq n$, let e_j be the j 'th standard basis vector in \mathbb{R}^n , and define $w_j = (u_1(j), u_2(j), \dots, u_k(j)) \in \mathbb{R}^k$, so $w_j(i) = u_i(j)$ for $1 \leq i \leq k$. For $1 \leq i \leq k$, let e'_i be the i 'th standard basis vector in \mathbb{R}^k . We will show that D is a V-cone by confirming that

$$D = \text{cone}(\{(e_j, w_j) : 1 \leq j \leq n\} \cup \{(-e_j, -w_j) : 1 \leq j \leq n\} \cup \{(0, e'_i) : 1 \leq i \leq k\}).$$

Call the convex cone on the right side D' . For $1 \leq i \leq k$ and $1 \leq j \leq n$, $u_i \bullet e_j = u_i(j) = w_j(i) \leq w_j(i)$ and $u_i \bullet (-e_j) = -u_i(j) = -w_j(i) \leq -w_j(i)$, so that $\pm(e_j, w_j) \in D$. For $1 \leq i, r \leq k$, $u_i \bullet 0 = 0 \leq e'_r(i)$ since $e'_r(i)$ is 0 or 1. These inequalities prove that all generators of D' lie in the convex cone D , so $D' \subseteq D$.

For the reverse inclusion, fix $(x, y) \in D$ satisfying the conditions in (11.2). For any real z , let $\text{sgn}(z) = 1$ if $z \geq 0$, $\text{sgn}(z) = -1$ if $z < 0$, so that $z = |z| \text{sgn}(z)$ in all cases. In Exercise 49, we ask the reader to prove the identity

$$(x, y) = \sum_{j=1}^n |x(j)|(\text{sgn}(x(j))e_j, \text{sgn}(x(j))w_j) + \sum_{i=1}^k (y(i) - u_i \bullet x)(0, e'_i). \quad (11.3)$$

This identity expresses (x, y) as a linear combination of generators of D' , where all the

scalars $|x(j)|$ and $y(i) - u_i \bullet x$ are nonnegative by the assumed conditions on (x, y) . It follows that $(x, y) \in D'$, completing the proof that $D = D'$ and D is a V-cone.

Comparing the original definitions of C and D , we see that $C \times \{0\}^k = D \cap (\mathbb{R}^n \times \{0\}^k)$. By invoking the intersection lemma in the form given at the end of §11.17, we conclude that $C \times \{0\}^k$ is a V-cone. Then C itself is also a V-cone.

11.19 Projection Lemma for H-Cones

One can readily check that the image of any V-cone in \mathbb{R}^n under any linear map is also a V-cone. To help prove that V-cones are the same as H-cones, we first need to prove that H-cones satisfy a special case of this *projection property*. Specifically, consider the linear projection map $p : \mathbb{R}^n \rightarrow \mathbb{R}^{n-1}$ given by $p(x) = (x(1), \dots, x(n-1))$ for $x \in \mathbb{R}^n$. Given any H-cone C in \mathbb{R}^n , we will prove that $p[C] = \{p(x) : x \in C\}$ is also an H-cone. Geometrically, $p[C]$ is the projection of C onto \mathbb{R}^{n-1} (identified with $\mathbb{R}^{n-1} \times \{0\}$ in \mathbb{R}^n).

By definition, there exist $k \in \mathbb{N}^+$ and $u_1, \dots, u_k \in \mathbb{R}^n$ with

$$C = \{x \in \mathbb{R}^n : u_i \bullet x \leq 0 \text{ for } 1 \leq i \leq k\}. \quad (11.4)$$

Moreover, $z \in \mathbb{R}^{n-1}$ lies in $p[C]$ iff there exists $r \in \mathbb{R}$ with $(z, r) \in C$. We need to pass from the system of linear inequalities (11.4) defining C to a new system of linear inequalities (involving only the first $n-1$ variables) with solution set $p[C]$. The idea is to take carefully chosen positive combinations of the original linear inequalities to *eliminate* the last variable x_n .

To implement this elimination process, write the index set $I = \{1, 2, \dots, k\}$ for the original list of inequalities as the disjoint union of subsets

$$I_0 = \{i \in I : u_i(n) = 0\}; \quad I_+ = \{i \in I : u_i(n) > 0\}; \quad I_- = \{i \in I : u_i(n) < 0\}.$$

For $i \in I_0$, let $u'_i = (u_i(1), u_i(2), \dots, u_i(n-1)) \in \mathbb{R}^{n-1}$. For $i \in I_+$ and $j \in I_-$, let $v_{ij} = u_i(n)u_j - u_j(n)u_i$. Each v_{ij} has n 'th component $u_i(n)u_j(n) - u_j(n)u_i(n) = 0$; let $v'_{ij} \in \mathbb{R}^{n-1}$ be v_{ij} with this zero deleted. Finally, let

$$D = \{z \in \mathbb{R}^{n-1} : u'_i \bullet z \leq 0 \text{ for } i \in I_0, \text{ and } v'_{ij} \bullet z \leq 0 \text{ for } i \in I_+, j \in I_-\}.$$

By definition, D is an H-cone in \mathbb{R}^{n-1} , so it will suffice to prove that $p[C] = D$.

To prove $p[C] \subseteq D$, we fix $z \in p[C]$ and check that $z \in D$. There is $r \in \mathbb{R}$ with $x = (z, r) \in C$. On one hand, for $i \in I_0$, we know $u_i \bullet x \leq 0$. Since $u_i(n) = 0$, the left side of this inequality does not involve r , and we get $u'_i \bullet z \leq 0$. On the other hand, fix $i \in I_+$ and $j \in I_-$. Multiply the known inequality $u_j \bullet x \leq 0$ by the scalar $u_i(n) \geq 0$, multiply the known inequality $u_i \bullet x \leq 0$ by the scalar $-u_j(n) \geq 0$, and add the resulting inequalities. We thereby see that $v_{ij} \bullet x \leq 0$. As noted above, $v_{ij}(n) = 0$ by construction, so we also have $v'_{ij} \bullet z \leq 0$. This means that z satisfies all the inequalities defining D , so that $z \in D$.

To prove $D \subseteq p[C]$, we fix $z \in D$ and must prove the existence of $r \in \mathbb{R}$ with $x = (z, r) \in C$. On one hand, for $i \in I_0$, the known inequality $u'_i \bullet z \leq 0$ guarantees the needed inequality $u_i \bullet x \leq 0$ for any choice of r , since $u_i(n) = 0$. On the other hand, for any fixed $i \in I_+$, the needed inequality $u_i \bullet x \leq 0$ will hold for $x = (z, r)$ iff $\sum_{s=1}^{n-1} u_i(s)z(s) + u_i(n)r \leq 0$ iff

$$r \leq u_i(n)^{-1} \sum_{s=1}^{n-1} (-u_i(s)z(s))$$

(recall $u_i(n) > 0$). Similarly, for any fixed $j \in I_-$, $u_j \bullet x \leq 0$ holds iff

$$r \geq u_j(n)^{-1} \sum_{s=1}^{n-1} (-u_j(s)z(s))$$

(note that the inequality is reversed since $u_j(n) < 0$). Combining these observations, we see that the remaining inequalities needed to ensure $x \in C$ will hold for $x = (z, r)$ iff

$$u_j(n)^{-1} \sum_{s=1}^{n-1} (-u_j(s)z(s)) \leq r \leq u_i(n)^{-1} \sum_{s=1}^{n-1} (-u_i(s)z(s))$$

for all $i \in I_+$ and $j \in I_-$. This collection of constraints on r is equivalent to the single condition

$$\max_{j \in I_-} \left[u_j(n)^{-1} \sum_{s=1}^{n-1} (-u_j(s)z(s)) \right] \leq r \leq \min_{i \in I_+} \left[u_i(n)^{-1} \sum_{s=1}^{n-1} (-u_i(s)z(s)) \right].$$

Now, there exists an $r \in \mathbb{R}$ satisfying this condition iff

$$\max_{j \in I_-} \left[u_j(n)^{-1} \sum_{s=1}^{n-1} (-u_j(s)z(s)) \right] \leq \min_{i \in I_+} \left[u_i(n)^{-1} \sum_{s=1}^{n-1} (-u_i(s)z(s)) \right],$$

which holds iff

$$u_j(n)^{-1} \sum_{s=1}^{n-1} (-u_j(s)z(s)) \leq u_i(n)^{-1} \sum_{s=1}^{n-1} (-u_i(s)z(s))$$

for all $i \in I_+$ and $j \in I_-$. Multiplying by $u_i(n)u_j(n) < 0$, the inequality just written is equivalent to

$$u_i(n) \sum_{s=1}^{n-1} (-u_j(s)z(s)) \geq u_j(n) \sum_{s=1}^{n-1} (-u_i(s)z(s)),$$

which holds iff

$$\sum_{s=1}^{n-1} (u_i(n)u_j(s)z(s) - u_j(n)u_i(s)z(s)) \leq 0$$

iff $(u_i(n)u_j - u_j(n)u_i) \bullet (z, r) \leq 0$ iff $v'_{ij} \bullet z \leq 0$. All of these last inequalities are true, since $z \in D$, and we see therefore that $z \in p[C]$. So $D \subseteq p[C]$.

We have now proved that for an H-cone C in \mathbb{R}^n , the projection $p[C]$ onto \mathbb{R}^{n-1} is an H-cone in \mathbb{R}^{n-1} . Iterating this result, we deduce that *for an H-cone C in \mathbb{R}^n and any $k \leq n$, the projection*

$$p_k[C] = \{z \in \mathbb{R}^{n-k} : \exists r_1, \dots, r_k \in \mathbb{R}, (z, r_1, \dots, r_k) \in C\} = \{z \in \mathbb{R}^{n-k} : \exists y \in \mathbb{R}^k, (z, y) \in C\}$$

is an H-cone in \mathbb{R}^{n-k} .

11.20 All V-Cones Are H-Cones

We now use the projection lemma for H-cones to prove that *every V-cone is an H-cone*. Given a V-cone $C \subseteq \mathbb{R}^n$, there are $k \in \mathbb{N}^+$ and $v_1, \dots, v_k \in \mathbb{R}^n$ with

$$C = \{c_1v_1 + \dots + c_kv_k : c_i \in \mathbb{R}, c_i \geq 0 \text{ for } 1 \leq i \leq k\}.$$

The key observation is that we can convert this V-cone into an H-cone in the higher-dimensional space \mathbb{R}^{n+k} . To do this, define

$$D = \left\{ (x, y) \in \mathbb{R}^{n+k} : x = \sum_{i=1}^k y(i)v_i \text{ and } y(i) \geq 0 \text{ for } 1 \leq i \leq k \right\}. \quad (11.5)$$

To see why D is an H-cone, let e_j and e'_i denote standard basis vectors in \mathbb{R}^n and \mathbb{R}^k , respectively. Define $u_i = (0, -e'_i)$ for $1 \leq i \leq k$ and $w_j = (e_j, (-v_1(j), \dots, -v_k(j)))$ for $1 \leq j \leq n$. One may check that (x, y) satisfies the conditions in the definition of D iff $u_i \bullet (x, y) \leq 0$ for $1 \leq i \leq k$ and $w_j \bullet (x, y) \leq 0$ for $1 \leq j \leq n$ and $(-w_j) \bullet (x, y) \leq 0$ for $1 \leq j \leq n$. So D is the solution set of a finite system of homogeneous linear inequalities, hence D is an H-cone.

To finish the proof, we need only observe that $x \in \mathbb{R}^n$ lies in C iff there exists $y \in \mathbb{R}^k$ with $(x, y) \in D$. By invoking the projection lemma in the form given at the end of §11.19, we conclude that C is an H-cone.

11.21 Finite Intersections of Closed Half-Spaces

Now that we know V-cones are the same as H-cones, we can prove the following theorem characterizing finite intersections of closed half-spaces. *A subset C of \mathbb{R}^n has the form $C = H_1 \cap H_2 \cap \dots \cap H_s$ for some $s \in \mathbb{N}$ and closed half-spaces H_i iff there exist $k, m \in \mathbb{N}$ and $v_1, \dots, v_k, w_1, \dots, w_m \in \mathbb{R}^n$ with $C = \text{conv}(\{v_1, \dots, v_k\}) + \text{cone}(\{w_1, \dots, w_m\})$.* The plus symbol here denotes the sum of sets in \mathbb{R}^n , namely $A + B = \{a + b : a \in A, b \in B\}$.

To begin the proof, first assume that

$$C = \text{conv}(\{v_1, \dots, v_k\}) + \text{cone}(\{w_1, \dots, w_m\}) \quad (k, m \in \mathbb{N}, v_i, w_j \in \mathbb{R}^n). \quad (11.6)$$

A typical point $v \in C$ looks like

$$v = c_1v_1 + \dots + c_kv_k + d_1w_1 + \dots + d_mw_m \quad (11.7)$$

where $c_i, d_j \geq 0$ are scalars such that $\sum_{i=1}^k c_i = 1$. Identifying \mathbb{R}^n with $\mathbb{R}^n \times \{0\}$, we can regard each v_i and w_j as a vector in \mathbb{R}^{n+1} . In \mathbb{R}^{n+1} , define a V-cone

$$D = \text{cone}(v_1 + e_{n+1}, \dots, v_k + e_{n+1}, w_1, \dots, w_m).$$

D consists of all points $v \in \mathbb{R}^{n+1}$ that can be written in the form

$$v = c_1v_1 + \dots + c_kv_k + d_1w_1 + \dots + d_mw_m + \left(\sum_{i=1}^k c_i \right) e_{n+1}, \quad (11.8)$$

for some scalars $c_i, d_j \geq 0$. Let H_0 denote the affine hyperplane $\mathbb{R}^n \times \{1\}$ in \mathbb{R}^{n+1} . We can obtain exactly those points in the intersection $D \cap H_0$ by choosing scalars $c_i, d_j \geq 0$ in (11.8) such that $\sum_{i=1}^k c_i = 1$. Since these are precisely the conditions imposed on the scalars in (11.7), we see that $D \cap H_0$ is the translate $C + e_{n+1}$.

As shown in §11.20, D is an H-cone, so D is a finite intersection of certain linear half-spaces H'_1, \dots, H'_s in \mathbb{R}^{n+1} . It follows that

$$C + e_{n+1} = D \cap H_0 = (H'_1 \cap H_0) \cap \dots \cap (H'_s \cap H_0).$$

Translating back to $\mathbb{R}^n \times \{0\}$, we see that $C = \bigcap_{i=1}^s ((H'_i \cap H_0) - e_{n+1})$, where each set in the intersection is readily seen to be a closed half-space in \mathbb{R}^n . So C has been expressed as a finite intersection of closed half-spaces.

Conversely, assume C is the intersection of finitely many closed half-spaces in \mathbb{R}^n , say

$$C = \{x \in \mathbb{R}^n : u_i \bullet x \leq b_i \text{ for } 1 \leq i \leq p\}$$

for some $p \in \mathbb{N}$ and $u_i \in \mathbb{R}^n$. In \mathbb{R}^{n+1} , define an H-cone

$$D = \{z \in \mathbb{R}^{n+1} : (u_i - b_i e_{n+1}) \bullet z \leq 0 \text{ for } 1 \leq i \leq p \text{ and } (-e_{n+1}) \bullet z \leq 0\}.$$

The last linear inequality in this definition amounts to requiring that the $(n+1)$ 'th coordinate of each point in D be nonnegative. Observe that for points $z \in \mathbb{R}^{n+1}$ of the form $(x, 1)$, the other linear inequalities in the definition of D hold iff $u_i \bullet x \leq b_i$ for each i . It follows from these remarks that $D \cap H_0 = C + e_{n+1}$, where (as before) H_0 denotes the affine hyperplane $\mathbb{R}^n \times \{1\}$ in \mathbb{R}^{n+1} .

As shown in §11.18, we can write $D = \text{cone}(\{z_1, \dots, z_s\})$ for some $s \in \mathbb{N}$ and $z_j \in \mathbb{R}^{n+1}$. By reordering the z_i 's and rescaling some z_i 's by positive scalars if needed, we can assume that $z_i(n+1) = 1$ for $1 \leq i \leq k$ and $z_i(n+1) = 0$ for $k < i \leq s$. Then a typical element of D is a positive combination

$$\sum_{i=1}^k c_i z_i + \sum_{i=k+1}^s d_i z_i \quad (c_i, d_i \geq 0).$$

If we want to find the points of $D \cap H_0$, we need to restrict the c_i 's and d_i 's to those nonnegative scalars that make the $(n+1)$ 'th coordinate equal to 1. By choice of the z_i 's, the restriction needed is precisely $\sum_{i=1}^k c_i = 1$. Since $D \cap H_0 = C + e_{n+1}$, we see that the points in C are exactly those points that can be written in the form

$$\sum_{i=1}^k c_i (z_i - e_{n+1}) + \sum_{i=k+1}^s d_i z_i \quad \left(c_i, d_i \geq 0, \sum_{i=1}^k c_i = 1 \right).$$

This means that

$$C = \text{conv}(\{z_1 - e_{n+1}, \dots, z_k - e_{n+1}\}) + \text{cone}(\{z_{k+1}, \dots, z_s\}),$$

where we view all generators as elements of \mathbb{R}^n . So C has been expressed in the form (11.6).

At last, we can prove the promised theorem that $C \subseteq \mathbb{R}^n$ is the convex hull of finitely many points iff C is a bounded intersection of finitely many closed half-spaces. Assume $C = \text{conv}(\{v_1, \dots, v_k\})$. This is the special case of (11.6) with $m = 0$, so C is an intersection of finitely many closed half-spaces. Moreover, any convex combination $w = \sum_{i=1}^k c_i v_i$ with $c_i \geq 0$ and $\sum_{i=1}^k c_i = 1$ will satisfy

$$\|w\| \leq \sum_{i=1}^k |c_i| \cdot \|v_i\| \leq \sum_{i=1}^k \|v_i\|,$$

so C is bounded. Conversely, assume C is a bounded set that is the intersection of finitely many closed half-spaces. We have proved that C can be written in the form (11.6), and we may assume no w_i is zero. If $m > 0$, then $\{dw_1 : d \geq 0\}$ would be an unbounded subset of C . So $m = 0$, and $C = \text{conv}(\{v_1, \dots, v_k\})$ as needed.

11.22 Convex Functions

We conclude the chapter with a brief introduction to convex functions. Let C be a convex subset of \mathbb{R}^n . A function $f : C \rightarrow \mathbb{R}$ is called *convex* iff $f(cx + (1 - c)y) \leq cf(x) + (1 - c)f(y)$ for all $x, y \in C$ and all $c \in [0, 1]$. (Compare to the definition of an affine map, where equality was required to hold for all $c \in \mathbb{R}$.) To see how this definition is related to the idea of convexity for sets, define the *epigraph* of f to be the set $\text{epi}(f) = \{(x, z) \in \mathbb{R}^{n+1} : x \in C, z \geq f(x)\}$, which consists of all points in \mathbb{R}^{n+1} “above” the graph of f . We claim f is a convex function iff $\text{epi}(f)$ is a convex set.

To verify this, first assume $\text{epi}(f)$ is a convex set. Given $x, y \in C$ and $c \in [0, 1]$, the points $(x, f(x))$ and $(y, f(y))$ are in $\text{epi}(f)$. By convexity of the epigraph, we know

$$c(x, f(x)) + (1 - c)(y, f(y)) = (cx + (1 - c)y, cf(x) + (1 - c)f(y)) \in \text{epi}(f).$$

The definition of the epigraph now gives $f(cx + (1 - c)y) \leq cf(x) + (1 - c)f(y)$, so f is a convex function. Conversely, assume f is a convex function. Fix $c \in [0, 1]$ and two points $(x, z), (y, w) \in \text{epi}(f)$, where $x, y \in \mathbb{R}^n$ and $z, w \in \mathbb{R}$. We know $z \geq f(x)$ and $w \geq f(y)$. By convexity of f ,

$$cz + (1 - c)w \geq cf(x) + (1 - c)f(y) \geq f(cx + (1 - c)y),$$

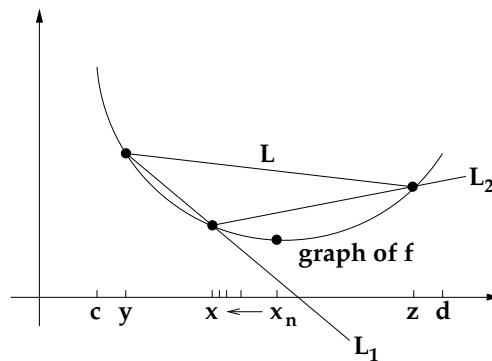
and hence $c(x, z) + (1 - c)(y, w) = (cx + (1 - c)y, cz + (1 - c)w) \in \text{epi}(f)$. So $\text{epi}(f)$ is a convex set. For $C \subseteq \mathbb{R}$, the convexity of f (or $\text{epi}(f)$) means that *for any real $a < b$ in C , the graph of f on the interval $[a, b]$ always lies weakly below the line segment joining $(a, f(a))$ to $(b, f(b))$.*

Convexity is a rather strong condition to impose on a function. For instance, *for all real $c < d$, a convex function $f : [c, d] \rightarrow \mathbb{R}$ must be continuous on the open interval (c, d)* . Figure 11.4 illustrates the proof. We verify continuity of f at an arbitrary point $x \in (c, d)$. Having fixed x , choose $y < x$ in (c, d) and $z > x$ in (c, d) . Draw the line L_1 of slope m_1 through $(y, f(y))$ and $(x, f(x))$ and the line L_2 of slope m_2 through $(x, f(x))$ and $(z, f(z))$. Since $(x, f(x))$ lies below the line L of slope m joining $(y, f(y))$ and $(z, f(z))$, it is geometrically evident that $m_1 \leq m \leq m_2$, so that L_1 is below L_2 to the right of x . Now consider any sequence of points x_n in the open interval (x, z) that converges to x . To see that f is right-continuous at x , we must show $\lim_{n \rightarrow \infty} f(x_n) = f(x)$. By using the remark at the end of the last paragraph with $a = x$ and $b = z$, we see that every point $(x_n, f(x_n))$ lies below the line L_2 . On the other hand, using the same remark with $a = y$ and $b = x_n$, we see that $(x, f(x))$ lies below the line joining $(y, f(y))$ to $(x_n, f(x_n))$, and hence $(x_n, f(x_n))$ lies above the line L_1 . Each of the lines L_1 and L_2 is the graph of a continuous (in fact, affine) function whose limit as x_n approaches x is $f(x)$. By the sandwich theorem for limits, x_n must also converge to x as n goes to infinity. Similar remarks prove the left-continuity of f at x . However, f need not be continuous at the endpoints c and d of its domain.

We have proved (§11.11) that convex sets are closed under convex combinations of their elements. Applying this remark to $\text{epi}(f)$ for a convex function $f : C \rightarrow \mathbb{R}$, we obtain a result called *Jensen's inequality*. Let $x_1, \dots, x_k \in C$ and $c_1, \dots, c_k \in \mathbb{R}$ with each $c_i \geq 0$ and $\sum_{i=1}^k c_i = 1$. Since $(x_i, f(x_i)) \in \text{epi}(f)$ for all i , the convex combination $(\sum_{i=1}^k c_i x_i, \sum_{i=1}^k c_i f(x_i))$ lies in $\text{epi}(f)$. Hence,

$$f\left(\sum_{i=1}^k c_i x_i\right) \leq \sum_{i=1}^k c_i f(x_i). \quad (11.9)$$

This is the *discrete form of Jensen's inequality*. By applying this inequality to Riemann

**FIGURE 11.4**

A Convex Function Must Be Continuous on an Open Interval.

sums approximating a Riemann integral, we can obtain a version of *Jensen's inequality for integrals*. Specifically, suppose $f : (c, d) \rightarrow \mathbb{R}$ is convex and $g : [0, 1] \rightarrow \mathbb{R}$ is an integrable function taking values in (c, d) . Then

$$f\left(\int_0^1 g(x) dx\right) \leq \int_0^1 f(g(x)) dx.$$

This follows from (11.9) by setting $c_i = 1/k$ and $x_i = g(i/k)$ for $1 \leq i \leq k$, which causes the two sums to approximate the two integrals just written, and then taking the limit of the inequality as k goes to infinity. By continuity of f , the limit of the left side of (11.9) is $f(\int_0^1 g(x) dx)$. Similarly, the right side approaches $\int_0^1 f(g(x)) dx$, where the integrability of $f \circ g$ follows since f is continuous. In the context of probability theory, sums involving convex combinations and integrals on $[0, 1]$ are special cases of expectations of random variables on a probability space. In this setting, Jensen's inequality becomes $f(E[X]) \leq E[f(X)]$, where f is convex, X is a random variable, and E denotes expected value.

11.23 Derivative Tests for Convex Functions

We have not yet given any examples of convex functions. Every affine function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is certainly convex; when $n = 1$, the graph of such a function is a line. The absolute value function ($f(x) = |x|$ for $x \in \mathbb{R}$) is convex, since it is visually apparent that the epigraph of f is convex. More generally, suppose $x_0 < x_1 < x_2 < \dots < x_n$ and y_0, y_1, \dots, y_n are given real numbers. Drawing the line segments from (x_{i-1}, y_{i-1}) to (x_i, y_i) for $1 \leq i \leq n$ will give the graph of a function $f : [x_0, x_n] \rightarrow \mathbb{R}$. This graph consists of a sequence of line segments having slopes $m_i = (y_i - y_{i-1})/(x_i - x_{i-1})$. One readily checks that f is convex iff $m_1 \leq m_2 \leq \dots \leq m_n$, i.e., the slopes of successive line segments are weakly increasing.

This suggests that a *differentiable* function $f : (c, d) \rightarrow \mathbb{R}$ is convex iff the first derivative $f' : (c, d) \rightarrow \mathbb{R}$ is an increasing function. To explain this geometrically, assume f is convex, and fix $y < x$ in (c, d) as shown in Figure 11.4. The secant line L_1 joining $(y, f(y))$ to $(x, f(x))$ has some slope m_1 . Convexity tells us that the graph of f between y and x is below this secant line, so it is geometrically evident that the tangent line to the graph of f at y has slope weakly less than m_1 . In symbols, $f'(y) \leq m_1$. But, as we argued earlier, the

graph of f to the right of x must lie above the secant line L_1 , so that the tangent line to the graph of f at x has slope weakly greater than m_1 . Thus, $m_1 \leq f'(x)$, so $f'(y) \leq f'(x)$ and f' is increasing.

Conversely, assume f' is increasing. To verify convexity of f , it suffices to show that for all $y < x < z$ in (c, d) , the point $(x, f(x))$ lies weakly below the line segment L joining $(y, f(y))$ and $(z, f(z))$ (see Figure 11.4). Let the secant lines L_1 and L_2 shown in the figure have respective slopes m_1 and m_2 ; then the point $(x, f(x))$ is weakly below L iff $m_1 \leq m_2$. By the mean value theorem, we can find a point y_1 with $y < y_1 < x$ and $f'(y_1) = m_1$; similarly, there is z_1 with $x < z_1 < z$ and $f'(z_1) = m_2$. Since $y_1 < z_1$ and f' is increasing, $m_1 \leq m_2$ as needed.

By a theorem from calculus, we deduce that a *twice-differentiable* function $f : (c, d) \rightarrow \mathbb{R}$ is convex iff the second derivative $f'' : (c, d) \rightarrow \mathbb{R}$ is nonnegative; i.e., $f''(x) \geq 0$ for all $x \in (c, d)$. This provides a convenient criterion for testing convexity of many functions. For example, $f(x) = e^x$ for $x \in \mathbb{R}$, $g(x) = x^2$ for $x \in \mathbb{R}$, and $h(x) = -\ln x$ for $x > 0$, are all convex functions by the second derivative test. In elementary calculus, these functions are often described as being *concave up*. For convex $C \subseteq \mathbb{R}^n$, we say that a function $f : C \rightarrow \mathbb{R}$ is *concave* iff $-f$ is convex iff $f(cx + (1 - c)y) \geq cf(x) + (1 - c)f(y)$ for all $x, y \in C$ and all $c \in [0, 1]$. When $n = 1$, our concave functions are called *concave down* in elementary calculus.

The second derivative test extends to multivariable functions as follows. Suppose C is an open convex subset of \mathbb{R}^n , and $f : C \rightarrow \mathbb{R}$ has continuous second-order partial derivatives. For each $z \in C$, define the *Hessian matrix* $H_z \in M_n(\mathbb{R})$ by letting $H_z(i, j) = f_{x_i x_j}(z) = \frac{\partial^2 f}{\partial x_i \partial x_j}(z)$ for $1 \leq i, j \leq n$. Then f is convex on C iff for all $z \in C$, H_z is a positive semidefinite matrix (which means $v \bullet H_z v = \sum_{i,j} f_{x_i x_j}(z)v_i v_j \geq 0$ for all $v \in \mathbb{R}^n$). The proof is outlined in Exercise 68; the idea is to reduce to the one-variable case by composing f with affine maps $g : (c, d) \rightarrow C$ and using the chain rule to compute $(f \circ g)''$.

11.24 Summary

Let F be a field in which $1_F + 1_F \neq 0_F$, and let V be an n -dimensional F -vector space.

1. *Types of Linear Combinations.* A *linear combination* of $v_1, \dots, v_k \in V$ is a vector of the form $c_1 v_1 + \dots + c_k v_k$ with all $c_i \in F$. This vector is: (a) an *affine* combination iff $\sum_{i=1}^k c_i = 1_F$; (b) a *convex* combination iff $F = \mathbb{R}$, all $c_i \geq 0$, and $\sum_{i=1}^k c_i = 1$; (c) a *positive* combination iff $F = \mathbb{R}$ and all $c_i \geq 0$.
2. *Definitions of Special Sets and Maps.* Table 11.1 summarizes the definitions of the structured sets and maps studied in this chapter.
3. *Characterizations of Linear Subspaces.* Every k -dimensional linear subspace of F^n is: (a) $\text{Sp}_F(S)$ for some linearly independent set S of size k ; (b) the null space of some $(n - k) \times n$ matrix; (c) the range of some $n \times k$ matrix; (d) the solution set of a system of $n - k$ homogeneous linear equations in n unknowns; (e) the intersection of $n - k$ linear hyperplanes in F^n .
4. *Characterizations of Affine Sets.* Every nonempty affine subset X of F^n of affine dimension k is: (a) $\text{aff}(S)$ for some affinely independent set S of size $k+1$; (b) $u+W$ for a unique k -dimensional linear subspace W (the direction subspace of X) and all $u \in X$; (c) the solution set of a system of $n - k$ (possibly non-homogeneous)

TABLE 11.1

Definitions of Affine and Convex Concepts.

Concept	Definition	Additional Properties
linear subspace W	$0_V \in W;$ $x, y \in W \Rightarrow x + y \in W;$ $x \in W, c \in F \Rightarrow cx \in W.$	W closed under linear combinations
affine set X	$x, y \in X, c \in F \Rightarrow$ $cx + (1 - c)y \in X.$	X closed under affine combinations; $X = u + W$ for dir. subspace W
convex set C	$x, y \in C, t \in [0, 1] \Rightarrow$ $tx + (1 - t)y \in C.$	C closed under convex combinations
cone C	$0 \in C;$ $x \in C, t \geq 0 \Rightarrow tx \in C$	C convex iff C closed under +
convex cone C	$0 \in C,$ $x, y \in C \Rightarrow x + y \in C,$ $x \in C, t \geq 0 \Rightarrow tx \in C.$	C closed under positive combinations
V-cone C	$C = \text{cone}(S)$ for some finite $S \subseteq \mathbb{R}^n$	same as H-cone
H-cone C	$C = \{x \in \mathbb{R}^n : u_i \bullet x \leq 0$ for $u_1, \dots, u_k \in \mathbb{R}^n\}$	same as V-cone
linear span $\text{Sp}_F(S)$	{lin. combs. of vectors in $S\}$	intersection of all subspaces $W \supseteq S$
affine span $\text{aff}(S)$	{aff. combs. of vectors in $S\}$	intersection of all affine sets $X \supseteq S$
convex hull $\text{conv}(S)$	{conv. combs. of points in $S\}$	intersection of all convex sets $C \supseteq S$
convex cone $\text{cone}(S)$	{pos. combs. of points in $S\}$	intersection of all convex cones $C \supseteq S$
linear hyperplane	$\{x \in \mathbb{R}^n : u \bullet x = 0\}$	2 choices for unit vector u
affine hyperplane	$\{x \in \mathbb{R}^n : u \bullet x = c\}$	unit vector u is unique when $c > 0$
closed half-space	$\{x \in \mathbb{R}^n : u \bullet x \leq c\}$	other half-space: $\{x \in \mathbb{R}^n : u \bullet x \geq c\}$
open half-space	$\{x \in \mathbb{R}^n : u \bullet x < c\}$	other half-space: $\{x \in \mathbb{R}^n : u \bullet x > c\}$
linear map T	$x, y \in V, c \in F \Rightarrow$ $T(x + y) = T(x) + T(y),$ $T(cx) = cT(x)$	T preserves linear combinations
affine map U	$x, y \in V, c \in F \Rightarrow$ $U(cx + (1 - c)y) =$ $cU(x) + (1 - c)U(y)$	U preserves affine combinations, $U(x) = T(x) + b$ (T linear)
convex function f (convex domain C)	$x, y \in C, t \in [0, 1] \Rightarrow$ $f(tx + (1 - t)y) \leq$ $tf(x) + (1 - t)f(y)$	continuous on open interval Jensen's inequality holds
concave function g (convex domain C)	$x, y \in C, t \in [0, 1] \Rightarrow$ $g(tx + (1 - t)y) \geq$ $tg(x) + (1 - t)g(y)$	continuous on open interval $-g$ is convex

linear equations in n unknowns; (d) the intersection of $n - k$ affine hyperplanes in F^n .

5. *Affine Independence, Bases, and Dimension.* A list $L = (v_0, v_1, \dots, v_k)$ of vectors in V is *affinely independent* iff for all $c_i \in F$ such that $\sum_{i=0}^k c_i = 0$, if $\sum_{i=0}^k c_i v_i = 0$ then all $c_i = 0$. L is affinely independent iff $(v_1 - v_0, \dots, v_k - v_0)$ is linearly independent. An affine basis of an affine set X is an affinely independent set whose affine span is X . The dimension of X is the dimension of its direction subspace, which is one less than the size of any affine basis of X . Affinely independent subsets of X can be extended to an affine basis; affine spanning sets

for X contain an affine basis; no affinely independent set is larger than an affine basis for X ; and any function defined on an affine basis of X extends uniquely to an affine map with domain X . Each point $v \in X$ has unique barycentric coordinates expressing v as an affine combination of a given ordered affine basis of X .

6. *Carathéodory's Theorem.* For all $S \subseteq \mathbb{R}^n$, every element of $C = \text{conv}(S)$ is a convex combination of at most $n + 1$ elements of S . One consequence is that convex hulls of closed and bounded (compact) sets are also closed and bounded.
7. *Separation by Hyperplanes.* If C is a closed, bounded, convex subset of \mathbb{R}^n and D is a closed, convex subset of \mathbb{R}^n disjoint from C , there is an affine hyperplane H in \mathbb{R}^n such that C and D lie in opposite open half-spaces of H . So, closed convex sets are the intersection of all open (or all closed) half-spaces containing them.
8. *Theorems on Generation vs. Intersection.* For all subsets C of \mathbb{R}^n :
 - (a) C has the form $\text{Sp}(S)$ for some set S iff C is an intersection of linear hyperplanes.
 - (b) C has the form $\text{aff}(S)$ for some set S iff C is an intersection of affine hyperplanes.
 - (c) C is closed and has the form $\text{conv}(S)$ for some set S iff C is an intersection of closed half-spaces.
 - (d) C has the form $\text{conv}(S)$ for some finite set S iff C is a bounded intersection of finitely many closed half-spaces.
 - (e) C has the form $\text{conv}(S) + \text{cone}(T)$ for some finite sets S and T iff C is an intersection of finitely many closed half-spaces.
 - (f) C has the form $\text{cone}(T)$ for some finite set T iff C is an intersection of finitely many linear half-spaces (i.e., V-cones are the same as H-cones).
9. *Intersection Formula for V-Cones.* Given a V-cone $C = \text{cone}(\{v_1, \dots, v_m\}) \subseteq \mathbb{R}^n$, the intersection $C \cap (\mathbb{R}^{n-1} \times \{0\})$ is the V-cone

$$D = \text{cone}(\{v_i : v_i(n) = 0\} \cup \{v_i(n)v_j - v_j(n)v_i : v_i(n) > 0, v_j(n) < 0\}).$$

10. *Projection Formula for H-Cones.* Given an H-cone

$$C = \{x \in \mathbb{R}^n : u_i \bullet x \leq 0 \text{ for } 1 \leq i \leq k\},$$

the projection $p[C] = \{z \in \mathbb{R}^{n-1} : \exists r \in \mathbb{R}, (z, r) \in C\}$ is the H-cone

$$D = \{z \in \mathbb{R}^{n-1} : u_i \bullet (z, 0) \leq 0 \text{ for } i \text{ with } u_i(n) = 0 \text{ and} \\ (u_i(n)u_j - u_j(n)u_i) \bullet (z, 0) \leq 0 \text{ for } i, j \text{ with } u_i(n) > 0, u_j(n) < 0\}.$$

11. *Theorems on Convex Functions.* A function f with convex domain is convex iff its epigraph $\text{epi}(f) = \{(x, y) : x \in C, y \geq f(x)\}$ is a convex set. For convex functions f and scalars $c_i \geq 0$ summing to 1, $f(\sum_i c_i x_i) \leq \sum_i c_i f(x_i)$ and $f(\int_0^1 g(x) dx) \leq \int_0^1 f(g(x)) dx$ whenever these expressions are defined (Jensen's inequality). A convex function whose domain is an open interval of \mathbb{R} must be continuous. For one-variable functions f such that f' exists, f is convex iff f' is increasing. If f'' exists, f is convex iff $f'' \geq 0$. If $C \subseteq \mathbb{R}^n$ is open and convex and $f : C \rightarrow \mathbb{R}$ has continuous second partial derivatives on C , f is convex iff the Hessian matrix $H_z = (f_{x_i, x_j}(z))_{1 \leq i, j \leq n}$ is positive semidefinite at all points of C .

11.25 Exercises

In these exercises, assume V and W are vector spaces over a field F with $\dim(V) = n$ unless otherwise stated.

1. (a) Prove that a subspace W of V is closed under linear combinations. (b) Prove that a subset W of V that is closed under linear combinations must be a subspace of V .
2. (a) Given $A \in M_{m,n}(F)$, prove that the range of A is a subspace of F^m . (b) Check that $T : F^n \rightarrow F^m$, given by $T(x) = Ax$ for $x \in F^n$, is an F -linear map. (c) Confirm that $\ker(T)$ is the null space of A , and $\text{img}(T)$ is the range of A .
3. (a) Prove that the intersection of any collection of linear subspaces of V is a linear subspace. (b) Prove that the intersection of any collection of affine subsets of V is affine. (c) Prove that the intersection of any collection of convex subsets of \mathbb{R}^n is convex. (d) Prove that the intersection of any collection of convex cones in \mathbb{R}^n is a convex cone.
4. Given subsets W_1, \dots, W_k of V , the *sum* of these subsets is $W_1 + \dots + W_k = \{x_1 + \dots + x_k : x_i \in W_i\}$. (a) Prove that the sum of linear subspaces is a linear subspace. (b) Prove that the sum of affine sets is affine. (c) Prove that the sum of convex sets is convex. (d) Prove that the sum of convex cones is a convex cone.
5. Let X and Y be subsets of V . (a) Prove or disprove: if X , Y , and $X \cup Y$ are all linear subspaces of V , then $X \subseteq Y$ or $Y \subseteq X$. (b) Prove or disprove: if X , Y , and $X \cup Y$ are all affine subsets of V , then $X \subseteq Y$ or $Y \subseteq X$. (c) Prove or disprove: if X , Y , and $X \cup Y$ are all convex ($V = \mathbb{R}^n$), then $X \subseteq Y$ or $Y \subseteq X$.
6. Given vector spaces V_1, \dots, V_k and subsets $S_i \subseteq V_i$, the *product* of these subsets is $S_1 \times \dots \times S_k = \{(x_1, \dots, x_k) : x_i \in S_i\}$. (a) Prove that the product of linear subspaces is a linear subspace. (b) Prove that the product of affine sets is an affine set. (c) Prove that the product of convex sets is a convex set. (d) Prove that the product of convex cones is a convex cone.
7. Let $S = \{v_1, \dots, v_k\} \subseteq V$. (a) Prove that $\text{Sp}(S)$ is a linear subspace of V . (b) Prove that $\text{Sp}(S)$ is the intersection of all subspaces of V that contain S .
8. Prove: (a) every m -dimensional subspace of F^n is the range of some matrix $A \in M_{n,m}(F)$; (b) every k -dimensional subspace of V is the kernel of some linear map $T : V \rightarrow F^{n-k}$; (c) every k -dimensional subspace of V is the image of some linear map $S : F^k \rightarrow V$.
9. Prove: for $A \in M_{m,n}(F)$ and $b \in F^m$, the set $\{x \in F^n : Ax = b\}$ is an affine subset of F^n .
10. Let S be a nonempty subset of V . (a) Prove S is a linear subspace of V iff $S + S = S$ and $cS = S$ for all nonzero $c \in F$. (b) Prove: if S is an affine set, then the direction subspace of S is $S + (-1)S$. (c) Prove: for all affine S and all $c \in F$, cS is an affine set.
11. We say that two nonempty affine sets $S, T \subseteq V$ are *parallel* iff S and T have the same direction subspace. (a) Show that two affine lines in \mathbb{R}^2 are parallel iff the lines are both vertical or both lines have the same slope. (b) Show that parallelism is an equivalence relation on the set of nonempty affine subsets of V , such that the equivalence class of S consists of all translates $u + S$ for $u \in V$.

12. Let X be an affine subset of V . (a) Prove: for all affinely independent $S \subseteq X$, there exists an affine basis of X containing S . (b) Prove: for all sets $T \subseteq X$ with $\text{aff}(T) = X$, there exists an affine basis of X contained in T .
13. (a) Show that any subset of an affinely independent set is also affinely independent. (b) Show that the maximum size of an affinely independent subset of the n -dimensional space V is $n + 1$.
14. Show that $\{v_0, v_1, \dots, v_k\}$ is affinely dependent iff some v_i is an affine combination of the v_j with $j \neq i$.
15. Given $S = \{v_0, v_1, \dots, v_k\} \subseteq \mathbb{R}^n$, let v'_k be an affine combination of v_0, \dots, v_k in which v_k appears with nonzero coefficient. Let S' be S with v_k replaced by v'_k . (a) Show $\text{aff}(S') = \text{aff}(S)$. (b) Show S' is affinely independent iff S is affinely independent.
16. Given $A \in M_{m,n}(F)$, let B be obtained from A by appending a column of ones to the right end of A . Show that the rows of A are affinely independent in F^n iff the rows of B are linearly independent in F^{n+1} .
17. *Points in General Position.* We say that points $v_1, \dots, v_k \in \mathbb{R}^n$ are *in general position* iff every subset of $\{v_1, \dots, v_k\}$ of size $n + 1$ or less is affinely independent. (a) For $n = 2$, show that distinct points $v_1, \dots, v_k \in \mathbb{R}^2$ are in general position iff no three v_i 's are collinear. (b) State and prove a result similar to (a) when $n = 3$. (c) Let c_1, \dots, c_k be distinct real numbers. Set $v_i = (c_i^n, c_i^{n-1}, \dots, c_i^2, c_i)$ for $1 \leq i \leq k$. Prove $v_1, \dots, v_k \in \mathbb{R}^n$ are in general position. [Hint: Use the fact that the Vandermonde matrix (see Exercise 30 in Chapter 5) has nonzero determinant.]
18. (a) Prove: for $n \in \mathbb{N}^+$, it is impossible to write \mathbb{R}^n as a union of finitely many affine hyperplanes. (b) Use (a) and induction on k to prove that for all $k, n \in \mathbb{N}^+$, there exist k points in general position in \mathbb{R}^n (see Exercise 17).
19. (a) Assume $1_F + 1_F \neq 0_F$ in the field F . Prove $S \subseteq V$ is affine iff $cx + (1 - c)y \in S$ for all $x, y \in S$ and all $c \in F$. (b) Give an example to show that the result in (a) can fail if $1 + 1 = 0$ in F . (c) Prove: for any field F , $S \subseteq V$ is affine iff $ax + by + cz \in S$ for all $x, y, z \in S$ and all $a, b, c \in F$ with $a + b + c = 1_F$.
20. Prove: for $A \in M_{m,n}(F)$ and $b \in F^m$, the map $U : F^n \rightarrow F^m$ given by $U(x) = Ax + b$ is an affine map. Show U is linear iff $b = 0$.
21. Prove: (a) the identity map on any affine set is an affine isomorphism; (b) the composition of two affine maps (resp. affine isomorphisms) is an affine map (resp. affine isomorphism); (c) the inverse of an affine isomorphism is also an affine isomorphism.
22. Let $U : X \rightarrow Y$ be an affine map between affine sets. (a) Prove $U[S]$ is affine for all affine $S \subseteq X$, and $U^{-1}[T]$ is affine for all affine $T \subseteq Y$. (b) Assume $X = \mathbb{R}^n$ and $Y = \mathbb{R}^m$. Prove $U[S]$ is convex for all convex $S \subseteq X$, and $U^{-1}[T]$ is convex for all convex $T \subseteq Y$.
23. Let $U : V \rightarrow W$ be an affine map. Prove: for all $S \subseteq V$, $U[\text{aff}(S)] = \text{aff}(U[S])$. Deduce as a special case that $u + \text{aff}(S) = \text{aff}(u + S)$ for all $u \in V$ and $S \subseteq V$.
24. (a) Prove: for every affine map $U : V \rightarrow W$ between vector spaces V and W , there exists a unique linear map $S : V \rightarrow W$ and a unique $b \in W$ with $U = T_b \circ S$, where $T_b : W \rightarrow W$ is translation by b . (b) Prove or disprove: for every affine map $U : V \rightarrow W$ between any vector spaces V and W , there exists a unique linear map $S : V \rightarrow W$ and a unique $c \in V$ with $U = S \circ T_c^{-1}$.

25. (a) Check that the map U defined at the end of §11.10 is affine and sends each v_i to y_i . (b) Prove the UMP for affine bases of X (see §11.10) by invoking the analogous UMP for a linear basis of the direction subspace of X (cf. (6.2)).
26. (a) Assume $1_F + 1_F \neq 0_F$ in the field F . Prove $U : V \rightarrow W$ is an affine map iff

$$U(cx + (1 - c)y) = cU(x) + (1 - c)U(y)$$

for all $x, y \in V$ and all $c \in F$. (b) Give an example to show that the result in (a) can fail if $1 + 1 = 0$ in F . (c) Prove: for any field F , $U : V \rightarrow W$ is an affine map iff $U(ax + by + cz) = aU(x) + bU(y) + cU(z)$ for all $x, y, z \in V$ and all $a, b, c \in F$ with $a + b + c = 1_F$.

27. Let X and Y be affine sets, let $C = \{v_0, \dots, v_k\} \subseteq X$, and let $U : X \rightarrow Y$ be an affine map. (a) Prove: if C is affinely independent and U is injective, then $U[C]$ is affinely independent. (b) Prove: if $X = \text{aff}(C)$ and U is surjective, then $Y = \text{aff}(U[C])$. (c) Prove: if C is an affine basis of X and U is bijective, then $U[C]$ is an affine basis of Y .
28. Let $B = (v_0, v_1, v_2) = ((1, 0, 1), (2, 3, 1), (1, 1, 2))$, and let $X = \text{aff}(B)$. (a) Prove B is an affinely independent ordered list. (b) Find the Cartesian coordinates of the point that has barycentric coordinates $(1/3, 1/2, 1/6)$ relative to B . (c) Find the barycentric coordinates relative to B of the point with Cartesian coordinates $(0, 0, 4)$.
29. Let $B = ((2, 0, 0), (1, 1, 0), (0, 1, 3), (-1, -1, 1))$, which is an affine basis of \mathbb{R}^3 . Given $v = (x, y, z) \in \mathbb{R}^3$, find the barycentric coordinates of v relative to B .
30. For $k \in \mathbb{N}$, let $\Delta = \text{conv}(\{e_1, e_2, \dots, e_{k+1}\}) \subseteq \mathbb{R}^{k+1}$. Show that Δ is a k -dimensional simplex such that Cartesian coordinates of points in Δ coincide with barycentric coordinates of points in Δ relative to the ordered affine basis (e_1, \dots, e_{k+1}) . Describe Δ as the solution set of a system of linear inequalities.
31. Show that for any two k -dimensional simplexes Δ_1 and Δ_2 in \mathbb{R}^n , there exists an affine isomorphism of \mathbb{R}^n mapping Δ_1 onto Δ_2 .
32. *Barycenter of a Simplex.* Let $\{v_0, \dots, v_k\}$ be affinely independent in \mathbb{R}^n . The *barycenter* of the simplex $\Delta = \text{conv}(\{v_0, \dots, v_k\})$ is the point in Δ all of whose barycentric coordinates (relative to the v_i 's) are $1/(k+1)$. (a) Show that when $k = 1$, the barycenter of Δ is the midpoint of the line segment with endpoints v_0 and v_1 . (b) Find the barycenter of the triangle $\text{conv}(\{(1, 2), (3, 5), (2, -4)\})$, and illustrate in a sketch. (c) Find the barycenter of $\text{conv}(\{0, e_1, \dots, e_n\})$ in \mathbb{R}^n . (d) For $0 \leq i \leq k$, let $\Delta_i = \text{conv}((\{v_0, \dots, v_k\} \sim \{v_i\}) \cup \{w\})$, where w is the barycenter of Δ . Prove that $\Delta = \bigcup_{i=0}^k \Delta_i$.
33. The *graph* of a function $T : V \rightarrow W$ is the set

$$G(T) = \{(x, T(x)) : x \in V\} \subseteq V \times W.$$

- (a) Show that the graph of a linear map is a linear subspace of $V \times W$. (b) Show that the graph of an affine map is an affine subset of $V \times W$. (c) Show that for every affine set S in F^n , there exist $m \leq n$, an affine map $U : F^{n-m} \rightarrow F^m$, and a linear isomorphism $P : F^n \rightarrow F^n$ that acts by permuting standard basis vectors, such that $P[S]$ is the graph of U .
34. (a) Show that for all linear hyperplanes H in \mathbb{R}^n , there exist exactly two unit vectors $u \in \mathbb{R}^n$ with $H = \{x \in \mathbb{R}^n : u \bullet x = 0\}$. (b) Show that for all affine hyperplanes H in \mathbb{R}^n not passing through zero, there is exactly one unit vector $u \in \mathbb{R}^n$ and one $c > 0$ with $H = \{x \in \mathbb{R}^n : u \bullet x = c\}$.

35. In \mathbb{R}^1 , explicitly describe all: (a) linear subspaces; (b) affine sets; (c) convex sets; (d) cones. Prove that the sets you describe, and no others, have each property.
36. (a) Prove that $C = \{(x, y, z) \in \mathbb{R}^3 : z^2 = x^2 + y^2\}$ is a cone. Is C a convex cone? Explain. (b) Prove $D = \{(x, y, z) \in \mathbb{R}^3 : z \geq \sqrt{x^2 + y^2}\}$ is a convex cone. (c) Is $E = \{(x, y, z) \in \mathbb{R}^3 : z^2 \geq x^2 + y^2\}$ a cone? Is E convex? Explain.
37. Fix $x_0 \in \mathbb{R}^n$ and $r \geq 0$. (a) Prove the *open ball* $B = \{x \in \mathbb{R}^n : d(x, x_0) < r\}$ is convex. (b) Is the *closed ball* $B' = \{x \in \mathbb{R}^n : d(x, x_0) \leq r\}$ convex? Explain. (c) Is the sphere $S = B' \setminus B$ convex? Explain.
38. Let S and T be convex subsets of \mathbb{R}^n . (a) Prove: for all $c \in \mathbb{R}$, cS is a convex set. (b) Prove: if S and T are convex sets and $c \in [0, 1]$, then $cS + (1 - c)T$ is a convex set. (c) Let $S = \text{conv}(\{(\cos(2\pi k/6), \sin(2\pi k/6)) : 0 \leq k < 6\})$ and $T = \text{conv}(\{(\cos(\pi/2 + 2\pi k/3), \sin(\pi/2 + 2\pi k/3)) : 0 \leq k < 3\})$. Sketch the sets $cS + (1 - c)T$ for $c \in \{0, 1/4, 1/2, 2/3, 1\}$.
39. Let $S \subseteq \mathbb{R}^n$ be a convex set. (a) Prove the closure of S (which can be defined as the set of $x \in \mathbb{R}^n$ with $d(x, S) = 0$) must be convex. (b) Prove the interior of S (the union of all open balls of \mathbb{R}^n contained in S) must be convex.
40. Let C be a nonempty compact (closed and bounded) subset of \mathbb{R}^n . (a) Given any nonempty set D in \mathbb{R}^n , define $f : C \rightarrow \mathbb{R}$ by $f(z) = d(z, D)$ for $z \in C$. Prove f is continuous, and conclude that there is $x \in C$ with $d(x, D) = d(C, D)$. (b) If D is compact, prove there is $x \in C$ and $y \in D$ with $d(x, y) = d(C, D)$. (c) Prove the conclusion of (b) holds for all closed D , even if D is unbounded. (d) Give an example where the conclusion of (b) fails for two closed, unbounded sets $C, D \subseteq \mathbb{R}^2$.
41. The separation theorem proved in §11.15 assumed that C was convex, closed, and bounded and that D was convex and closed. Give examples to show that the omission of any one of these five assumptions may cause the conclusion of the theorem to fail.
42. Complete the proof of the separation theorem in §11.15 by showing that $D \subseteq \{v \in \mathbb{R}^n : u \bullet v \geq c_2\}$ and indicating what adjustments are needed in the case $c_1 > c_2$.
43. Give an example of a convex subset of \mathbb{R}^2 that cannot be written as an intersection of any family of half-spaces (even allowing a mixture of open and closed half-spaces).
44. (a) Prove that the intersection of any family of closed subsets of \mathbb{R}^n is closed. (b) Prove that every closed half-space is a closed set. (c) Prove that every affine hyperplane is a closed set. (d) Prove that every affine set is a closed set. (e) Prove that every simplex is a closed set. (f) Prove that for fixed $x_0 \in \mathbb{R}^n$ and $r > 0$, the “closed ball” $\{x \in \mathbb{R}^n : d(x, x_0) \leq r\}$ is a closed set.
45. (a) Prove that if $S \subseteq \mathbb{R}^n$ is closed and bounded, then $\text{conv}(S)$ is a closed convex subset of \mathbb{R}^n . (Use Carathéodory’s theorem and sequential compactness of S and $[0, 1]$.) (b) Give an example of a countably infinite set S such that $\text{conv}(S)$ is not closed. (c) If S is bounded but not closed, must $\text{conv}(S)$ be bounded? Prove or give a counterexample.
46. *Radon’s Theorem.* Suppose $S = \{v_0, v_1, \dots, v_{n+1}\}$ is a set of $n+2$ vectors in \mathbb{R}^n . Prove there exist disjoint subsets T and U of S with $\text{conv}(T) \cap \text{conv}(U) \neq \emptyset$.
47. Fix $n, k \in \mathbb{N}$. For any convex sets $C, D \subseteq \mathbb{R}^{n+k}$, let $C +_{n,k} D$ be the set
- $$\{(x, y) \in \mathbb{R}^n \times \mathbb{R}^k : \exists u \in \mathbb{R}^k, \exists v \in \mathbb{R}^k, (x, u) \in C, (x, v) \in D, \text{ and } y = u + v\}.$$

- (a) Prove $C +_{n,k} D$ is convex. (b) Prove $+_{n,k}$ is an associative, commutative binary operation on the set of convex subsets of \mathbb{R}^{n+k} .
48. (a) Let $C, D \subseteq \mathbb{R}^n$ be convex sets. Define $C \# D = \bigcup_{t \in [0,1]} [(1-t)C \cap tD]$. Prove $C \# D$ is convex. [Hint: Study $\text{cone}(C + e_{n+1}) +_{n,1} \text{cone}(D + e_{n+1})$.] (b) Prove: if C and D are convex cones, then $C \# D = C \cap D$. (c) Prove: if C and D are convex cones, then $C + D = \text{conv}(C \cup D)$. (d) Give examples to show that (b) and (c) can fail if C and D are convex sets that are not cones.
49. This exercise checks some details of the proof in §11.18. (a) Show that D defined in (11.2) is a convex cone. (b) Confirm identity (11.3). (c) Confirm that $C \times \{0\}^k = D \cap (\mathbb{R}^n \times \{0\}^k)$. (d) Show that since $C \times \{0\}^k$ is a V-cone, C is a V-cone.
50. For the set D defined in (11.5), confirm that $(x, y) \in D$ iff (x, y) satisfies the linear inequalities given in §11.20.
51. Let $S = \{(0, 1, 2), (0, -1, 1), (-1, 0, 2), (3, 0, 0)\} \subseteq \mathbb{R}^3$. (a) Sketch S , $\text{Sp}(S)$, $\text{aff}(S)$, $\text{conv}(S)$, and $\text{cone}(S)$. (b) Give a linear inequality defining the closed half-space of $\text{aff}(S)$ containing the origin. (c) Give a description of $\text{cone}(S)$ as an H-cone.
52. Prove: for all $S \subseteq \mathbb{R}^n$, $\text{conv}(S) \subseteq \text{aff}(S) \cap \text{cone}(S)$. Does equality always hold?
53. (a) Prove the intersection of an H-cone in \mathbb{R}^n with any linear subspace of \mathbb{R}^n is an H-cone. (b) Prove the image of a V-cone in \mathbb{R}^n under a linear map is a V-cone. (c) Can you give a direct proof of (a) for V-cones, or (b) for H-cones, without using the theorem that V-cones and H-cones are the same?
54. Call a subset P of \mathbb{R}^n a *V-polyhedron* iff $P = \text{conv}(S)$ for a finite set $S \subseteq \mathbb{R}^n$. (a) Prove that the intersection of finitely many V-polyhedra is a V-polyhedron. (b) Prove that the intersection of a V-polyhedron with an affine subset of \mathbb{R}^n is a V-polyhedron.
55. Call a subset P of \mathbb{R}^n an *H-polyhedron* iff $P = H_1 \cap \dots \cap H_s$ for finitely many closed half-spaces H_i . (a) Prove that the image of an H-polyhedron under an affine map is an H-polyhedron. (b) Prove that the sum of finitely many H-polyhedra is an H-polyhedron.
56. Let $S = \{(-2, 1), (-1, 4), (0, -1), (1, 1), (2, 3), (3, 0)\} \subseteq \mathbb{R}^2$. (a) Sketch $\text{conv}(S)$ in \mathbb{R}^2 and $\text{cone}(S + e_3)$ in \mathbb{R}^3 . (b) Express $\text{cone}(S + e_3)$ as a specific intersection of linear half-spaces. (c) Express $\text{conv}(S)$ as a specific intersection of closed half-planes.
57. Let S be the solution set in \mathbb{R}^2 of the system of linear inequalities
- $$y - x \leq 3, \quad x \leq 2, \quad x + y \leq 3, \quad x - 2y \leq 4, \quad -y \leq 2, \quad 2x + y \leq -2, \quad -3x + y \leq 3.$$
- Find a finite set T with $S = \text{conv}(T)$.
58. Let S be the solution set in \mathbb{R}^3 of the linear inequalities $0 \leq z \leq 2$, $0 \leq y \leq z$, $y \leq x \leq (y+6)/3$, $x - y + z \leq 3$. Sketch S and find a finite set T with $S = \text{conv}(T)$.
59. Let $S = \{(\pm 3, 0, 0), (1, \pm 1, 0), (-1, \pm 1, 0), (\pm 1/2, 0, 1), (0, \pm 1/2, 1)\}$. Sketch $T = \text{conv}(S)$ and find a specific system of linear inequalities with solution set T .
60. Let C be the V-cone $\text{cone}(v_1, \dots, v_5)$ in \mathbb{R}^4 , where

$$v_1 = (3, 1, 1, 1), \quad v_2 = (0, 1, 2, 2), \quad v_3 = (1, 0, 2, 0),$$

$$v_4 = (-1, -1, 1, -1), \quad v_5 = (2, 2, 0, 3).$$

- (a) Follow the proof in §11.17 to find a finite set S with $C \cap (\mathbb{R}^3 \times \{0\}) = \text{cone}(S)$.
 (b) Find a finite set T with $C \cap (\mathbb{R}^2 \times \{0\}^2) = \text{cone}(T)$.
61. Let C be the H-cone in $\mathbb{R}^4 = \{(w, x, y, z) : w, x, y, z \in \mathbb{R}\}$ defined by the system of inequalities

$$2w - x + y \leq 0, \quad w + x + 2y + 2z \leq 0, \quad x - 3z \leq 0, \quad w + y - z \leq 0, \quad 2x + y + z \leq 0.$$

(a) Follow the proof in §11.19 to find a system of inequalities whose solution set is the projection of C onto $\mathbb{R}^3 \times \{0\}$. (b) Find a system of inequalities defining the projection of C onto $\mathbb{R}^2 \times \{0\}^2$. What is this projection?

62. Let C be the V-cone $\text{cone}(v_1, \dots, v_5)$ in \mathbb{R}^3 , where

$$v_1 = (1, 2, 1), \quad v_2 = (1, 4, 0), \quad v_3 = (4, 1, 0), \quad v_4 = (4, 3, -1), \quad v_5 = (3, 4, -1).$$

Sketch C and find an explicit description of C as an H-cone.

63. Let C be the H-cone defined by the inequalities

$$x - 3y - z \leq 0, \quad x + y - z \leq 0, \quad -11x + 7y - z \leq 0, \quad -x - 3y - z \leq 0.$$

Sketch C and find an explicit description of C as a V-cone.

64. In Figure 11.4, let the lines L , L_1 , and L_2 have respective slopes m , m_1 , and m_2 .
 (a) Carefully prove that $(x, f(x))$ lies weakly below L iff $m_1 \leq m \leq m_2$. (b) Use the definition to prove that for convex $C \subseteq \mathbb{R}$, $f : C \rightarrow \mathbb{R}$ is convex iff for all $y < x < z$ in C , $(x, f(x))$ lies weakly below the line segment joining $(y, f(y))$ and $(z, f(z))$.
65. Determine whether each function below is convex, concave, or neither.
 (a) $f(x) = -2x + 5$ for $x \in \mathbb{R}$; (b) $f(x) = x^3 - x$ for $x \in \mathbb{R}$; (c) $f(x) = \sin^2 x$ for $x \in \mathbb{R}$; (d) $f(x) = \cos x$ for $x \in [\pi/2, 3\pi/2]$; (e) $f(x) = \arctan(x)$ for $x \geq 0$; (f) $f(x) = \sqrt{1-x^2}$ for $x \in [-1, 1]$; (g) $f : [0, \infty) \rightarrow \mathbb{R}$ given by $f(x) = x^r$, for fixed $r \geq 1$; (h) $f : [0, \infty) \rightarrow \mathbb{R}$ given by $f(x) = x^r$, for fixed $r \in [0, 1)$; (i) $f : (0, \infty) \rightarrow \mathbb{R}$ given by $f(x) = x^r$, for fixed $r < 0$.
66. (a) Draw a diagram similar to Figure 11.4 and use it to prove that the convex function f is left-continuous at $x \in (c, d)$. (b) Give an example to show that a convex function $f : [c, d] \rightarrow \mathbb{R}$ may not be right-continuous at c or left-continuous at d .
67. Suppose $a < b$ in \mathbb{R} , $f : (c, d) \rightarrow \mathbb{R}$ is convex, and $g : [a, b] \rightarrow (c, d)$ is integrable.
 (a) Prove $(b-a)f(\int_a^b g(x) dx) \leq \int_a^b f((b-a)g(x)) dx$. (b) Give an example to show that (a) can fail if we omit $(b-a)$ from both sides of the inequality.
68. Let C be an open convex subset of \mathbb{R}^n , and let $f : C \rightarrow \mathbb{R}$ be a function with continuous second-order partial derivatives. (a) Show f is convex on C iff for all $c < d$ in \mathbb{R} and all affine maps $g : (c, d) \rightarrow C$, $f \circ g : (c, d) \rightarrow \mathbb{R}$ is a convex function. (b) Let $g : (c, d) \rightarrow C$ have the formula $g(t) = x + tv$ for $t \in (c, d)$, where $x, v \in \mathbb{R}^n$ are fixed. Show that $f \circ g$ is convex iff $v \bullet H_y v \geq 0$ for all y in the image of g . (c) Use (a) and (b) to prove that f is convex iff H_y is positive semidefinite for all $y \in C$.
69. (a) Prove $g : \mathbb{R}^n \rightarrow \mathbb{R}$, given by $g(x) = \|x\| = \sqrt{x \bullet x}$ for $x \in \mathbb{R}^n$, is convex.
 (b) Let $C = \{x \in \mathbb{R}^n : x(i) > 0 \text{ for } 1 \leq i \leq n\}$. Prove $f : C \rightarrow \mathbb{R}^n$, given by $f(x) = -(x(1)x(2)\cdots x(n))^{1/n}$ for $x \in C$, is convex.

70. Let $f : C \rightarrow \mathbb{R}$ be convex. (a) Prove: for all $r \in \mathbb{R}$, $\{x \in C : f(x) < r\}$ and $\{x \in C : f(x) \leq r\}$ are convex sets. (b) Must the sets $\{x \in C : f(x) > r\}$ and $\{x \in C : f(x) \geq r\}$ be convex? Prove or give a counterexample.
71. Prove: for all $k \geq 1$ and all real $z_1, \dots, z_k > 0$, $(z_1 + z_2 + \dots + z_k)/k \geq \sqrt[k]{z_1 z_2 \cdots z_k}$. (Apply Jensen's inequality to a certain convex function.)
72. Let B and C be convex sets, let $f : C \rightarrow \mathbb{R}$ and $g : \mathbb{R} \rightarrow \mathbb{R}$ be convex functions, and let $T : B \rightarrow C$ be a linear map. (a) Prove: if g is increasing, then $g \circ f$ is convex. (b) Give an example to show that $g \circ f$ may not be convex if g is not increasing. (c) Prove: $f \circ T$ is convex. (d) If T is only an affine map, must (c) be true? Explain.
73. (a) Suppose $f_1, \dots, f_k : C \rightarrow \mathbb{R}$ are convex functions. Prove every positive linear combination of these functions is convex. (b) Suppose $\{f_i : i \in I\}$ is a family of convex functions with domain C such that $g(x) = \sup_{i \in I} f_i(x)$ is finite for all $x \in C$. Prove $g : C \rightarrow \mathbb{R}$ is convex.
74. Fix real numbers $p, q > 1$ with $p^{-1} + q^{-1} = 1$. Prove: For $r, s \in \mathbb{R}^+$,

$$r^{1/p} s^{1/q} \leq r/p + s/q.$$

This page intentionally left blank

12

Ruler and Compass Constructions

Ancient Greek geometers spent a lot of time investigating certain construction problems that arise in plane geometry. Here are some examples of geometric construction problems.

- *Altitudes:* Given a line ℓ and a point P , construct the line through P that meets ℓ at a right angle.
- *Parallels:* Given a line ℓ and a point P not on ℓ , construct the line through P that is parallel to ℓ .
- *Polygons:* Given a circle with center O and a point P on the circle, construct a regular n -sided polygon inscribed in the circle and having P as one vertex. For example, when $n = 3, 4, 5, 6, 7$, we are trying to construct equilateral triangles, squares, regular pentagons, regular hexagons, and regular heptagons.
- *Angle Bisection:* Given two lines ℓ_1, ℓ_2 that meet at a point P at an angle θ , draw a line through P that makes an angle $\theta/2$ with ℓ_1 .
- *Angle Trisection:* Given two lines ℓ_1, ℓ_2 that meet at a point P at an angle θ , draw a line through P that makes an angle $\theta/3$ with ℓ_1 .
- *Cube Duplication:* Given a line segment \overline{PQ} , construct a line segment \overline{AB} such that a cube with \overline{AB} as one side has twice the volume of a cube with \overline{PQ} as one side.
- *Squaring the Circle:* Given a line segment \overline{PQ} , construct a line segment \overline{AB} such that the square with \overline{AB} as one side has the same area as the circle with \overline{PQ} as radius.

To solve these construction problems, the Greeks allowed the use of only two rudimentary tools: a *ruler*, which can draw the straight line passing through any two points; and a *compass*, which can draw the circle with a given center passing through a given point. The ruler *cannot* be used to measure distances, and the compass *cannot* be used to draw several circles with the same radius in different parts of the diagram. To emphasize these restrictions, the ruler is sometimes called a *straightedge*, and the compass is sometimes called a *collapsing* compass.

Perhaps surprisingly, a great variety of construction problems can be solved using only a ruler and compass. For example, the constructions of altitudes, parallels, angle bisectors, and regular n -gons for $n = 3, 4, 5, 6, 8$ are all possible. On the other hand, some of the other problems on the preceding list turn out to be *impossible* to solve using only a ruler and compass. The most famous of these problems are squaring the circle, trisecting an arbitrary angle, and duplicating the cube. It is also impossible to construct a regular heptagon (7-sided polygon) and many other regular polygons. One of the amazing achievements of Gauss was the discovery that a regular 17-sided polygon *can* be constructed using a ruler and compass.

This chapter develops the mathematical tools needed to prove the unsolvability or solvability of various ruler and compass constructions. Remarkably, deciding the solvability of a given geometric construction can ultimately be translated into a linear algebra question involving the dimensions of certain vector spaces. To see how this arises, we must first link

the geometric operations occurring in ruler and compass constructions to various arithmetic operations, notably the extraction of square roots. We then use field theory to characterize the numbers obtainable by means of these arithmetic operations. Any field K with subfield F can be regarded as a vector space with scalars coming from F . The fact that each such vector space has a unique *dimension* will turn out to be the key to proving the impossibility or possibility of our geometric constructions.

To implement this agenda, we begin by giving rigorous definitions of three different kinds of “constructible” numbers — one based on geometry, another based on arithmetic, and a third based on field theory. The core theorem of this subject asserts that the three notions of constructibility are all equivalent. Once we prove this theorem, we use it to analyze some of the famous construction problems mentioned above.

The main prerequisites for reading this chapter are the definitions of fields and vector spaces (see Chapter 1), basic facts about the dimension of vector spaces (see §1.8), an acquaintance with elementary Euclidean geometry and analytic geometry, and some knowledge of irreducible polynomials (Chapter 3).

12.1 Geometric Constructibility

We first give an informal description of ruler and compass constructions, then translate this description into a rigorous definition. Our geometric constructions occur in a two-dimensional Euclidean plane. We can use an x, y -coordinate system to identify this plane with \mathbb{R}^2 , the set of ordered pairs (a, b) with $a, b \in \mathbb{R}$. Alternatively, we can identify the plane with the set \mathbb{C} of complex numbers by letting the point (a, b) correspond to the complex number $a + ib$. The initial data for our geometric constructions consist of two points A and B in the plane corresponding to the complex numbers 0 and 1. Thus, A has coordinates $(0, 0)$, and B has coordinates $(1, 0)$.

We are now allowed to construct new points, lines, and circles in \mathbb{R}^2 by applying a finite sequence of actions from the following list.

- If two different points P and Q have already been constructed, we can draw the unique line through P and Q , which we denote by $L(P, Q)$.
- If two different points P and Q have already been constructed, we can draw the unique circle with center P passing through Q , which we denote by $C(P; Q)$.
- If two unequal, non-parallel lines ℓ_1 and ℓ_2 have been constructed, we can locate the unique point where these two lines intersect.
- If a line ℓ and a circle C have been constructed, we can locate all intersection points of this line and this circle (if any).
- If two unequal circles C_1 and C_2 have been constructed, we can locate all intersection points of these two circles (if any).

For our proofs, we need a more formal definition of geometrically constructible points. We recursively define a set GC of complex numbers, called the *geometrically constructible* numbers, by the following rules.

$$G0. \quad 0 \in GC \text{ and } 1 \in GC.$$

$$G1. \quad \text{If } P, Q, R, S \in GC, L(P, Q) \text{ and } L(R, S) \text{ are unequal lines, and } L(P, Q) \cap L(R, S) = \{T\}, \text{ then } T \in GC.$$

- G2. If $P, Q, R, S \in \text{GC}$ and $T \in L(P, Q) \cap C(R; S)$, then $T \in \text{GC}$.
- G3. If $P, Q, R, S \in \text{GC}$, $C(P; Q)$ and $C(R; S)$ are unequal circles, and $T \in C(P; Q) \cap C(R; S)$, then $T \in \text{GC}$.
- G4. The only numbers in GC are those that can be obtained by applying rules G0, G1, G2, and G3 a finite number of times.

We can also rephrase this recursive definition in the following iterative fashion. A complex number Q is in GC iff there is a finite sequence of points

$$P_0 = 0, P_1 = 1, P_2, P_3, \dots, P_k = Q \quad (k \geq 0)$$

such that, for all i with $2 \leq i \leq k$, there exist $r, s, t, u < i$ for which P_i is in $L(P_r, P_s) \cap L(P_t, P_u)$ or in $L(P_r, P_s) \cap C(P_t; P_u)$ or in $C(P_r; P_s) \cap C(P_t; P_u)$. (In the first and third alternatives, we require that the two lines or circles do not coincide.) We call the sequence P_0, P_1, \dots, P_k a *geometric construction sequence* for Q .

For example, the sequence $0, 1, -1, i\sqrt{3}, \sqrt{3}$ is a geometric construction sequence for $\sqrt{3}$, so $\sqrt{3} \in \text{GC}$. To verify this, first note that -1 is a point in the intersection of the line $L(0, 1)$ (the x -axis) and the circle $C(0; 1)$ (the unit circle). Next, $i\sqrt{3}$ is one of the two intersection points of $C(1; -1)$ and $C(-1; 1)$, since these circles have radius 2 and the point $(0, \sqrt{3})$ is a distance 2 away from both $(-1, 0)$ and $(1, 0)$. Finally, $\sqrt{3}$ is in $C(0; i\sqrt{3}) \cap L(0, 1)$.

12.2 Arithmetic Constructibility

Next we consider a notion of constructibility involving arithmetic operations. We want to study complex numbers that can be built up from the integers by performing the familiar arithmetic operations (addition, subtraction, multiplication, and division in \mathbb{C}) together with the operation of extracting square roots. Such numbers arise naturally, for example, in connection with the quadratic formula. Recall that for $a, b, c \in \mathbb{C}$ with a nonzero, the roots of the equation $ax^2 + bx + c$ in \mathbb{C} are

$$r_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a}, \quad r_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}.$$

We see that these roots can be built from the initial data a, b, c by performing arithmetic operations and taking square roots.

To formalize this idea, we recursively define a set AC of complex numbers, called the *arithmetically constructible* numbers, by the following rules.

- A0. $0 \in \text{AC}$ and $1 \in \text{AC}$.
- A1. If $a, b \in \text{AC}$ then $a + b \in \text{AC}$.
- A2. If $a \in \text{AC}$ then $-a \in \text{AC}$.
- A3. If $a, b \in \text{AC}$ then $a \cdot b \in \text{AC}$.
- A4. If $a \in \text{AC}$ is nonzero, then $a^{-1} \in \text{AC}$.
- A5. If $a \in \text{AC}$ and $b \in \mathbb{C}$ and $b^2 = a$, then $b \in \text{AC}$.
- A6. The only numbers in AC are those that can be obtained by applying rules A0 through A5 a finite number of times.

A less formal phrasing of rule A5 says that if $a \in AC$ then $\sqrt{a} \in AC$. Since every nonzero complex number has two complex square roots, the notation \sqrt{a} used here is ambiguous. However, this is not a serious difficulty, since rule A2 guarantees that both square roots of a must lie in AC .

As before, we can give an iterative version of the definition of AC , as follows. A complex number z is in AC iff there is a finite sequence of numbers

$$x_0 = 0, x_1 = 1, x_2, x_3, \dots, x_k = z \quad (k \geq 0)$$

such that, for all i with $2 \leq i \leq k$, either there exists $r < i$ with $x_i = -x_r$ or $x_i = x_r^{-1}$ ($x_r \neq 0$) or $x_i = \sqrt{x_r}$, or there exist $r, s < i$ with $x_i = x_r + x_s$ or $x_i = x_r \cdot x_s$. The sequence x_0, x_1, \dots, x_k is called an *arithmetic construction sequence* for z .

For example, an arithmetic construction sequence for $i\sqrt{3}$ is $0, 1, 2, 3, -3, i\sqrt{3}$, hence $i\sqrt{3} \in AC$. To build this sequence, we invoked rule A0 twice, then rule A1 twice, then rule A2, then rule A5. For a more elaborate illustration of arithmetic construction sequences, suppose $a, b, c \in AC$ with $a \neq 0$. Then the roots r_1, r_2 of the quadratic equation $ax^2 + bx + c$ (given above) are also in AC . For instance, we can see that $r_1 \in AC$ by producing an arithmetic construction sequence as follows. Begin the sequence with $0, 1, 2, 3, 4$, followed by the concatenation of arithmetic construction sequences for a, b , and c . Continue this sequence as follows:

$$b^2, ac, 4ac, -4ac, b^2 - 4ac, \sqrt{b^2 - 4ac}, -b, -b + \sqrt{b^2 - 4ac}, 2a, r_1.$$

One sees readily that each new term does arise from one or two previous terms by invoking one of the rules.

12.3 Preliminaries on Field Extensions

To continue, we need some facts about field extensions. Let K be a field (see §1.2 for the definition of a field). Recall (§1.4) that a *subfield* of K is a subset F of K such that $0_K \in F$, $1_K \in F$, and for all $a, b \in F$, $a + b \in F$ and $-a \in F$ and $a \cdot b \in F$ and (for $a \neq 0_K$) $a^{-1} \in F$. (The similarity of these closure conditions to rules A0 through A4 above is not accidental, and will be exploited below.) A subfield F of K becomes a field by restricting the sum and product operations on K to the subset F . We say K is an *extension field* of F iff F is a subfield of K . A *chain of fields* is a sequence $F_0 \subseteq F_1 \subseteq F_2 \subseteq \dots \subseteq F_n$ in which F_n is a field, and each F_i (for $0 \leq i < n$) is a subfield of F_n . For most of this chapter, the fields under consideration will be subfields of \mathbb{C} , so that the operations in these fields are ordinary addition and multiplication of complex numbers.

Now let K be any field with subfield F . A key observation is that we can regard K as a vector space over the field F , by taking vector addition to be the field addition $+ : K \times K \rightarrow K$, and taking scalar multiplication $s : F \times K \rightarrow K$ to be the restriction of the field multiplication $\cdot : K \times K \rightarrow K$ to the domain $F \times K$. In other words, we multiply a “scalar” $c \in F$ by a “vector” $v \in K$ by forming the product $c \cdot v$ in the field K . The vector space axioms follow by comparing the axioms for the field K (and subfield F) to the axioms for an F -vector space. Like any other F -vector space, K has a unique dimension viewed as a vector space over F . This dimension is denoted $[K : F]$ and is called the *degree of K over F* . For example, $K = \mathbb{C}$ is a two-dimensional vector space over its subfield $F = \mathbb{R}$ (since $(1, i)$ is an ordered basis of the real vector space \mathbb{C}), and so the degree $[\mathbb{C} : \mathbb{R}]$ is 2. On the other hand, $[\mathbb{R} : \mathbb{Q}] = \infty$, since \mathbb{R} is uncountable but every finite-dimensional \mathbb{Q} -vector

space is countable. (To be precise, $[\mathbb{R} : \mathbb{Q}]$ should be a certain infinite cardinal giving the \mathbb{Q} -dimension of \mathbb{R} . However, for our purposes in this chapter, we will not need this extra precision. Thus, the notation $[K : F] = \infty$ merely indicates that K is an infinite-dimensional F -vector space.)

Suppose $F \subseteq K \subseteq E$ is a chain of three fields. We will prove the *degree formula*

$$[E : F] = [E : K][K : F],$$

which will play a critical role in all subsequent developments. We assume in our proof that $[E : K] = m < \infty$ and $[K : F] = n < \infty$; if these finiteness assumptions are not met, one can show (Exercise 15) that both sides of the degree formula are ∞ . Since $[E : K] = m$, there exists an ordered basis $B = (y_1, y_2, \dots, y_m)$ for the K -vector space E ; we call B a *K -basis* of E to emphasize that scalars here come from K . Similarly, $[K : F] = n$ means there is an ordered F -basis $C = (z_1, \dots, z_n)$ for the F -vector space K . Consider the list of all products $z_j y_i$:

$$D = (z_1 y_1, z_1 y_2, \dots, z_1 y_m, z_2 y_1, \dots, z_2 y_m, \dots, z_n y_1, \dots, z_n y_m).$$

This is a list of mn elements, which are not yet known to be distinct. We claim D is an ordered F -basis of E , which will prove that $[E : F] = mn$.

First we show D spans E as an F -vector space. Given any $w \in E$, we first write $w = \sum_{i=1}^m k_i y_i$ for certain scalars $k_i \in K$, which is possible since B is a K -basis of E . Since each $k_i \in K$, and since C is an F -basis of K , we can find further scalars $f_{ij} \in F$ such that $k_i = \sum_{j=1}^n f_{ij} z_j$ for $1 \leq i \leq m$. Inserting these expressions into the formula for w and simplifying, we obtain

$$w = \sum_{i=1}^m \left(\sum_{j=1}^n f_{ij} z_j \right) y_i = \sum_{i=1}^m \sum_{j=1}^n f_{ij} (z_j y_i).$$

We have expressed an arbitrary vector $w \in E$ as an F -linear combination of the elements of the list D , so this list does span E as claimed.

Next, we prove that the list D is F -linearly independent. Assume $c_{ij} \in F$ are scalars for which $\sum_{i=1}^m \sum_{j=1}^n c_{ij} (z_j y_i) = 0$; we must prove all c_{ij} 's are zero. Using the distributive laws to regroup terms, our assumption can be written

$$0 = \sum_{i=1}^m \left(\sum_{j=1}^n c_{ij} z_j \right) y_i.$$

Each term in parentheses lies in the field K , since $z_j \in K$ and $c_{ij} \in F \subseteq K$. The list $B = (y_1, \dots, y_m)$ is known to be K -linearly independent, so we deduce that $\sum_{j=1}^n c_{ij} z_j = 0$ for all i with $1 \leq i \leq m$. Now we can invoke the known F -linear independence of the list $C = (z_1, \dots, z_n)$ to see that $c_{ij} = 0$ for all i, j with $1 \leq i \leq m$ and $1 \leq j \leq n$.

This completes the proof of the degree formula for a chain of three fields. Now consider a chain of fields of arbitrary finite length, say

$$F_0 \subseteq F_1 \subseteq F_2 \subseteq F_3 \subseteq \cdots \subseteq F_n.$$

Iteration of the previous result gives the *general degree formula*

$$[F_n : F_0] = [F_n : F_{n-1}] \cdots [F_3 : F_2][F_2 : F_1][F_1 : F_0] = \prod_{i=1}^n [F_i : F_{i-1}].$$

The last preliminary concept we need is the notion of a field extension generated by one element. Suppose K is a field with subfield F , and let $z \in K$. Define $F(z)$, the *field extension of F obtained by adjoining z* , to be the intersection of all subfields L of K such that $F \subseteq L$ and $z \in L$. It is routine to check that $F(z)$ is a subfield of K with $F \subseteq F(z)$ and $z \in F(z)$; and for any subfield M of K with $F \subseteq M$ and $z \in M$, $F(z) \subseteq M$. One can also verify that $F(z)$ consists of all elements of K that can be written in the form $f(z)g(z)^{-1}$, for some polynomials f, g with coefficients in F such that $g(z) \neq 0$. In this chapter, we will mainly be interested in the case where $K = \mathbb{C}$ and $F = \mathbb{Q}$, so $\mathbb{Q}(z)$ is the smallest subfield of the complex numbers containing the given complex number z .

12.4 Field-Theoretic Constructibility

Our final notion of constructibility involves field extensions of degree 2. We say that a complex number z has a *square root tower* iff there is a chain of fields

$$\mathbb{Q} = F_0 \subseteq F_1 \subseteq \cdots \subseteq F_n \subseteq \mathbb{C}$$

such that $z \in F_n$ and $[F_k : F_{k-1}] = 2$ for $1 \leq k \leq n$. Let SQC denote the set of all $z \in \mathbb{C}$ such that z has a square root tower. We observe that

$$z \in \text{SQC} \Rightarrow [\mathbb{Q}(z) : \mathbb{Q}] = 2^e \text{ for some } e \geq 0. \quad (12.1)$$

To prove this, first use the degree formula to see that $[F_n : \mathbb{Q}] = 2^n$. Since $z \in F_n$ and $\mathbb{Q} \subseteq F_n$, we have a chain of fields $\mathbb{Q} \subseteq \mathbb{Q}(z) \subseteq F_n$. It follows from the degree formula that $[\mathbb{Q}(z) : \mathbb{Q}]$ divides $[F_n : \mathbb{Q}] = 2^n$, so that $[\mathbb{Q}(z) : \mathbb{Q}]$ must be a power of 2. We warn the reader that the converse of (12.1) is *not* true (although this text will not develop the machinery needed to disprove the converse). However, the contrapositive of (12.1) is certainly true:

$$\text{if } [\mathbb{Q}(z) : \mathbb{Q}] \text{ is not a power of 2, then } z \notin \text{SQC}. \quad (12.2)$$

The main theorem of this chapter states that

$$\text{GC} = \text{AC} = \text{SQC}.$$

In other words, a number z is geometrically constructible iff z is arithmetically constructible iff z has a square root tower. Once this theorem is known, we can use the criterion (12.2) to prove the unsolvability of various geometric construction problems. (In fact, for such applications, it suffices to know the weaker result $\text{GC} \subseteq \text{AC} \subseteq \text{SQC}$.) Similarly, we can demonstrate the solvability of certain geometric constructions by building appropriate square root towers. This agenda will be carried out in §12.9 and §12.10, after we prove the main result $\text{GC} = \text{AC} = \text{SQC}$ in the next few sections.

12.5 Proof that $\text{GC} \subseteq \text{AC}$

In this section, we give the formal proof that $\text{GC} \subseteq \text{AC}$. The details are a bit tedious, but the intuition underlying the proof is simple: when computing intersections of lines and circles via analytic geometry, the new coordinates can always be computed from coordinates

of previous points by arithmetic operations and extractions of square roots. This is the key feature of the messy formulas derived below.

Our proof requires the following facts from plane analytic geometry.

- AG1. For every line ℓ in the plane, there exist real numbers a, b, s such that

$$\ell = \{(x, y) \in \mathbb{R}^2 : ax + by = s\};$$

here a and b are not both zero. If $\ell = L(P, Q)$ where $P = (x_0, y_0)$ and $Q = (x_1, y_1)$, then we may take $a = y_1 - y_0$, $b = x_0 - x_1$, and $s = ax_0 + by_0 = ax_1 + by_1$. Note especially that if $x_0, y_0, x_1, y_1 \in \text{AC}$, then we can choose a, b, s so that $a, b, s \in \text{AC}$.

- AG2. If C is a circle in \mathbb{R}^2 with center (c, d) and radius $r > 0$, then

$$C = \{(x, y) \in \mathbb{R}^2 : (x - c)^2 + (y - d)^2 = r^2\}.$$

If (e, f) is any point on C , then $r^2 = (e - c)^2 + (f - d)^2$. Note especially that if $c, d, e, f \in \text{AC}$, then $r^2, r \in \text{AC}$.

- AG3. Consider lines ℓ_1 and ℓ_2 with equations $ax + by = r$ and $cx + dy = s$. These lines are non-parallel (and, in particular, non-equal) iff $ad - bc \neq 0$. In this case, the unique intersection (x_0, y_0) of ℓ_1 and ℓ_2 is given by Cramer's rule (see §5.12):

$$x_0 = \frac{rd - bs}{ad - bc}, \quad y_0 = \frac{as - rc}{ad - bc}.$$

Note especially that if $a, b, c, d, r, s \in \text{AC}$, then $x_0, y_0 \in \text{AC}$.

- AG4. Suppose a line ℓ and a circle C have respective equations $ax + by = s$ and $(x - c)^2 + (y - d)^2 = r^2$. A point (x, y) lies on both ℓ and C iff both equations hold simultaneously. Suppose that b is nonzero. Substituting $y = (s - ax)/b$ into the equation of the circle, we see that x must satisfy the equation $(x - c)^2 + ((s - ax)/b - d)^2 = r^2$. This equation can be rewritten as $Ax^2 + Bx + D = 0$, where

$$A = 1 + \frac{a^2}{b^2}, \quad B = -2c + \frac{2ad}{b} - \frac{2as}{b^2}, \quad D = c^2 + d^2 - r^2 - \frac{2ds}{b} + \frac{s^2}{b^2}.$$

So x is a root of a linear or quadratic equation with the indicated coefficients, and $y = (s - ax)/b$ is determined by x . Note especially that if $a, b, s, c, d, r \in \text{AC}$, then $A, B, D \in \text{AC}$ and hence $x, y \in \text{AC}$. A similar analysis holds when a is nonzero.

- AG5. Suppose C_1 and C_2 are circles with equations $(x - a)^2 + (y - b)^2 = r^2$ and $(x - c)^2 + (y - d)^2 = s^2$. A point (x, y) lies on both circles iff both equations hold simultaneously. Subtracting the second equation from the first, we obtain a third equation

$$(2c - 2a)x + (2d - 2b)y = c^2 + d^2 + r^2 - a^2 - b^2 - s^2.$$

Assuming that the circles are non-concentric, this is the equation of a line ℓ . The first and second equations have the same solution set as the second and third equations, as is readily seen. It follows that the intersection points of C_1 and C_2 are the same as the intersection points of ℓ and C_2 . So we can find the coordinates of these points using the formulas from the preceding item. Note especially that if $a, b, c, d, r, s \in \text{AC}$, then the coordinates x and y of the intersection point(s) lie in AC .

Now we are ready to prove that $\text{GC} \subseteq \text{AC}$. Given $z = (x, y) = x + iy \in \text{GC}$, we know there is a geometric construction sequence z_0, z_1, \dots, z_k with last term z . Write $z_r = x_r + iy_r$ with $x_r, y_r \in \mathbb{R}$ for $0 \leq r \leq k$. We will prove that $x = x_k$ and $y = y_k$ are in AC by strong induction on k . Since $i \in \text{AC}$, it will follow that $z = x + iy$ is in AC . If $k = 0$, then $z_0 = 0 = 0 + 0i$, and 0 is in AC . If $k = 1$, then $z_1 = 1 = 1 + 0i$, and $0, 1 \in \text{AC}$. Let $k > 1$; we can assume inductively that for all $r < k$, $x_r \in \text{AC}$ and $y_r \in \text{AC}$. Now consider the three possible cases for z_k .

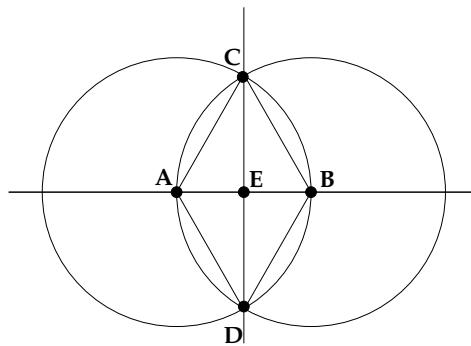
- *Case G1:* There exist $r, s, t, u < k$ such that $z_k = (x_k, y_k)$ is the unique point in the intersection of the line ℓ_1 through (x_r, y_r) and (x_s, y_s) and the line ℓ_2 through (x_t, y_t) and (x_u, y_u) . Since x_r, y_r , etc., all lie in AC by induction, we see that x_k and y_k lie in AC by AG1 and AG3.
 - *Case G2:* There exist $r, s, t, u < k$ such that $z_k = (x_k, y_k)$ is in the intersection of the line ℓ through (x_r, y_r) and (x_s, y_s) and the circle C with center (x_t, y_t) passing through (x_u, y_u) . Since x_r, y_r , etc., all lie in AC by induction, we see that x_k and y_k lie in AC by AG1, AG2, and AG4.
 - *Case G3:* There exist $r, s, t, u < k$ such that $z_k = (x_k, y_k)$ is in the intersection of the circle C_1 through (x_s, y_s) with center (x_r, y_r) and the circle C_2 through (x_u, y_u) with center (x_t, y_t) . Since the coordinates of these points all lie in AC by induction, we see that x_k and y_k lie in AC by AG2 and AG5.
-

12.6 Proof that $\text{AC} \subseteq \text{GC}$

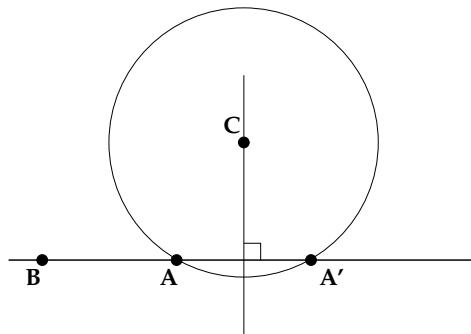
We have proved that GC is contained in AC by showing that the coordinates of the intersections of lines and circles can always be found from previous coordinates by arithmetic operations and square root extractions. To prove the reverse containment, we must show conversely that arithmetic operations and square root extractions can be performed geometrically with a ruler and compass. We present the necessary geometric constructions now, followed by the formal proof that $\text{AC} \subseteq \text{GC}$.

- CR1. *Given two points A and B , we can construct an equilateral triangle with side \overline{AB} , the midpoint of \overline{AB} , and the perpendicular bisector of \overline{AB} .* First draw the line connecting A and B ; then draw the circle centered at A passing through B ; then draw the circle centered at B passing through A . The two circles must meet¹ in two points, say C and D . Then ΔABC and ΔABD are equilateral triangles, \overleftrightarrow{CD} is the perpendicular bisector of \overleftrightarrow{AB} , and the intersection of \overleftrightarrow{CD} and \overleftrightarrow{AB} is the midpoint of \overleftrightarrow{AB} .

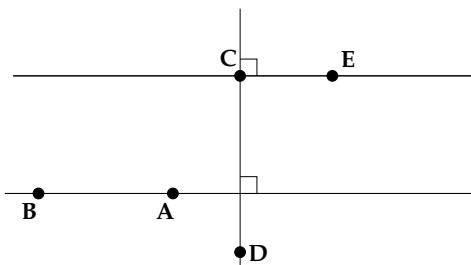
¹Like Euclid, we omit the proofs that the claimed intersection points really do exist.



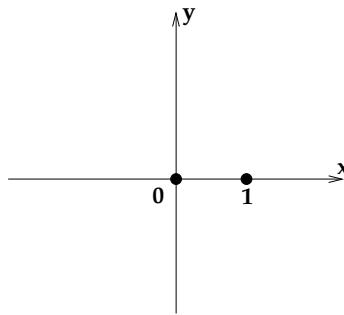
- CR2. Given three points A, B, C with $A \neq B$, we can construct the unique altitude through C perpendicular to \overleftrightarrow{AB} . Choose notation so that $C \neq A$. Draw \overleftrightarrow{AB} and the circle with center C passing through A . If this circle meets \overleftrightarrow{AB} at a second point $A' \neq A$, then the required altitude is the perpendicular bisector of $\overline{AA'}$, which can be built by CR1. If, instead, the circle is tangent to \overleftrightarrow{AB} at A , the required altitude is \overleftrightarrow{AC} . The construction works whether or not C is on the line \overleftrightarrow{AB} .



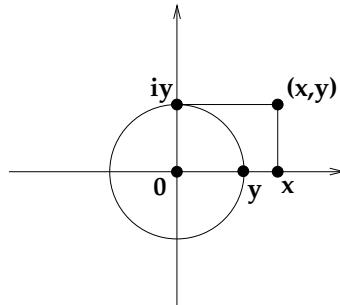
- CR3. Given three points A, B, C with $A \neq B$, we can construct the unique line through C parallel to \overleftrightarrow{AB} . Apply CR2 once to get a line \overleftrightarrow{CD} perpendicular to \overleftrightarrow{AB} . Apply CR2 again to get a line \overleftrightarrow{CE} perpendicular to \overleftrightarrow{CD} , hence parallel to \overleftrightarrow{AB} .



- CR4. Given the initial points 0 and 1, we can draw the x -axis and the y -axis. The x -axis is the line through 0 and 1. Now use CR2 to draw the y -axis.

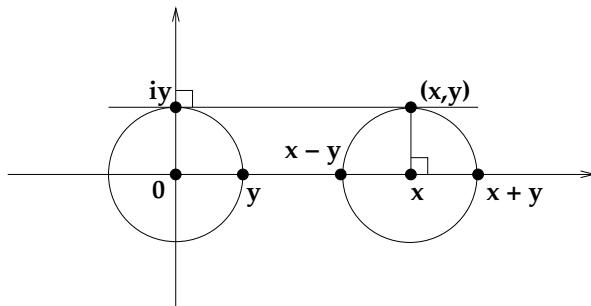


- CR5. For all real x, y , we have $x + iy = (x, y) \in \text{GC}$ iff $x = (x, 0) \in \text{GC}$ and $y = (y, 0) \in \text{GC}$. For if $x + iy = (x, y) \in \text{GC}$, we can use CR2 to drop altitudes to the x -axis and the y -axis to find the points $(x, 0)$ and $(0, y)$. Drawing the circle centered at 0 and passing through $(0, y)$ allows us to locate $(y, 0)$. The converse statement is proved similarly by reversing these steps.



The next few constructions show how to implement real arithmetic using a ruler and compass. We can then handle complex arithmetic by using CR5 to reduce to the real case.

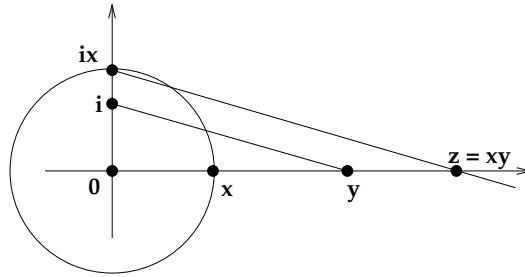
- CR6. For all real $x, y \in \text{GC}$, $x + y \in \text{GC}$ and $x - y \in \text{GC}$. Given x, y , locate the point $(x, y) = x + iy$ as in CR5. The circle with center x passing through $x + iy$ meets the real axis at the points $x + y$ and $x - y$. Hence these are in GC.



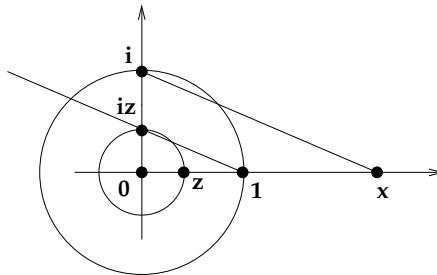
- CR7. For all real $y \in \text{GC}$, $-y \in \text{GC}$. Take $x = 0$ in CR6.

- CR8. For all real $x, y \in \text{GC}$, $xy \in \text{GC}$. Drawing circles with center 0 through 1 and x , we can locate i and ix on the imaginary axis. Draw the line through i and y , and then use CR3 to draw the line through ix parallel to this line. Let z be the point where this line meets the real axis. By considering the two similar right triangles in the figure, we see that $y/1 = z/x$, so that $z = xy$. Our figure assumes that x

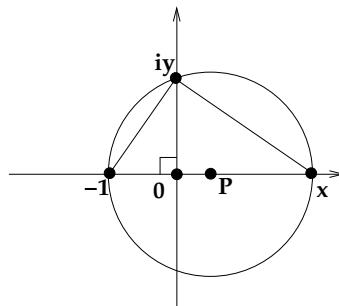
and y are positive. The reader can adapt this construction to the case of negative variables (alternatively, this can be deduced from the positive case using CR7).



- CR9. For all nonzero real $x \in \text{GC}$, $x^{-1} \in \text{GC}$. Draw i and the line through x and i . Then draw the line through 1 parallel to this line, which meets the imaginary axis in some point iz . Use a circle centered at 0 to find the real point z . Comparing similar right triangles again, we see that $z/1 = 1/x$, so $z = x^{-1}$ is constructible.



- CR10. For all positive real $x \in \text{GC}$, $\sqrt{x} \in \text{GC}$. Draw -1 , and then construct the midpoint P of the line segment joining -1 and x . Draw the circle through x with center P , and let iy be the point where this circle meets the imaginary axis. One can then construct y , which we claim is the positive square root of x . To see this, draw line segments from iy to -1 and from iy to x . The angle formed by these segments at iy is a right angle, since it is inscribed in a semicircle. It readily follows that the right triangle with vertices 0 , -1 , and iy is similar to the right triangle with vertices 0 , iy , and x . So $1/y = y/x$, which shows that $y^2 = x$ and $y = \sqrt{x}$.



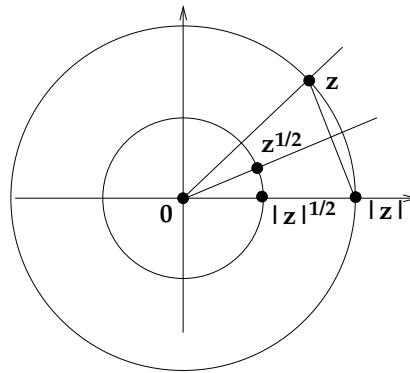
- CR11. For all complex $u, v \in \text{GC}$, we have $u + v \in \text{GC}$, $-u \in \text{GC}$, $uv \in \text{GC}$, and (if $u \neq 0$) $u^{-1} \in \text{GC}$. Write $u = a + ib$ and $v = c + id$ where $a, b, c, d \in \mathbb{R}$. By CR5,

the real numbers a, b, c, d are all in GC. Now,

$$\begin{aligned} u + v &= (a + c) + i(b + d), & -u &= (-a) + i(-b), \\ uv &= (ac - bd) + i(ad + bc), & u^{-1} &= \left(\frac{a}{a^2 + b^2} \right) + i \left(\frac{-b}{a^2 + b^2} \right). \end{aligned}$$

Using CR6, CR7, CR8, and CR9, we see that the real and imaginary parts of all these expressions are in GC. Hence, another application of CR5 shows that $u + v$, $-u$, uv , and u^{-1} are all in GC.

- CR12. For all complex $z \in \text{GC}$, $\sqrt{z} \in \text{GC}$. Write $z = re^{i\theta}$ in polar form, where $r \geq 0$ and θ are real numbers. Drawing the circle through z with center 0, we see that $r = |z| \in \text{GC}$. By CR10, we can construct the real square root of r . By CR1, we can construct the perpendicular bisector of the line segment joining z and r . Then the points w where this bisector intersects the circle with center 0 and radius \sqrt{r} are in GC. One may check that $w = \pm\sqrt{re^{i\theta/2}}$ in polar form, and these are the two complex square roots of z . Thus both square roots lie in GC.



We now have the tools needed to prove that $\text{AC} \subseteq \text{GC}$. Given $z \in \text{AC}$, we know there is an arithmetic construction sequence z_0, z_1, \dots, z_k with last term z . We proceed by strong induction on k . If $k = 0$, then $z = 0$ is in GC. If $k = 1$, then $z = 1$ is in GC. If $k > 1$, we can assume by induction that z_0, \dots, z_{k-1} are all in GC. Now if $z_k = z_r + z_s$ for some $r, s < k$, then $z_k \in \text{GC}$ follows by CR11. The same conclusion holds if $z_k = -z_r$ or $z_k = z_r z_s$ or $z_k = 1/z_r$ where $r, s < k$. Finally, if $z_k = \sqrt{z_r}$ for some $r < k$, then $z_k \in \text{GC}$ by CR12. Thus, in all cases, $z = z_k$ is in GC.

12.7 Algebraic Elements and Minimal Polynomials

Before proving that $\text{AC} = \text{SQC}$, we need some more facts about field extensions of the form $F(z)$. Suppose K is a field, F is a subfield of K , and $z \in K$. We say that z is *algebraic over F* iff there exists a nonzero polynomial $g \in F[x]$ satisfying $g(z) = 0_K$. We will prove that *if z is algebraic over F , then there exists a unique monic, irreducible polynomial $m \in F[x]$ such that $m(z) = 0_K$; and no polynomial in $F[x]$ of smaller degree than m has z as a root*. The polynomial m is called the *minimal polynomial of z over F* .

To begin the proof, assume $z \in K$ is algebraic over F . Since there exist nonzero polynomials in $F[x]$ having z as a root, we can pick a nonzero polynomial $m \in F[x]$

of minimum possible degree for which $m(z) = 0$. Dividing by the leading coefficient if needed, we can assume that m is monic. We claim m must be an irreducible polynomial in $F[x]$, i.e., there cannot exist any factorization $m = fh$ with $f, h \in F[x]$ each having lower degree than m . For, if such a factorization did exist, evaluating both sides at z would give $0_K = m(z) = f(z)h(z)$. Since K is a field, this forces $f(z) = 0_K$ or $h(z) = 0_K$; but either possibility contradicts minimality of the degree of m .

Now we prove the uniqueness of m . To get a contradiction, assume $p \neq m$ is another monic, irreducible polynomial in $F[x]$ with $p(z) = 0_K$. Since p and m are non-associate irreducible polynomials, $\gcd(p, m) = 1$. It follows (see Chapter 3) that $pr + ms = 1$ for some $r, s \in F[x]$. Evaluating both sides at z and recalling that $p(z) = 0_K = m(z)$, we obtain $0_K = 1_K$. But this contradicts the definition of a field.

We can use minimal polynomials to give an explicit basis for the field extension $F(z)$, viewed as an F -vector space, in the case where z is algebraic over F . We will prove: *if K is a field with subfield F , and $z \in K$ is algebraic over F with minimal polynomial $m \in F[x]$ of degree d , then $B = (1, z, z^2, \dots, z^{d-1})$ is an ordered F -basis for $F(z)$, and hence $[F(z) : F] = d = \deg(m)$.* Let W be the subspace of the F -vector space K spanned by the vectors in the list B . Note that W can also be described as the set of elements of the form $g(z)$, where $g \in F[x]$ is either zero or has degree less than d . We first show that B is an F -linearly independent list. Suppose $c_0, c_1, \dots, c_{d-1} \in F$ satisfy $c_0 + c_1z + \dots + c_{d-1}z^{d-1} = 0_K$. Then the polynomial $g = c_0 + c_1x + \dots + c_{d-1}x^{d-1}$ is in $F[x]$ and has z as a root. Since no polynomial of degree less than $d = \deg(m)$ has z as a root, it must be that g is the zero polynomial, so $c_0 = c_1 = \dots = c_{d-1} = 0$. We now know that B is a basis for the F -vector space W .

To finish the proof, we will show that $W = F(z)$. On one hand, since $F(z)$ is a subfield of K containing z and F , closure under sums and products shows that $F(z)$ must contain all powers of z and all F -linear combinations of these powers. In particular, $W \subseteq F(z)$. To prove the reverse inclusion $F(z) \subseteq W$, it suffices to show that W is a subfield of K with $F \subseteq W$ and $z \in W$. Keeping in mind that W is the set of F -linear combinations of the powers of z in B , it is routine to check that W contains F (hence contains 0_K and 1_K), that $z \in W$, and that W is closed under addition, additive inverses, and scalar multiplication by elements of F .

Is W closed under multiplication? Consider two arbitrary elements $u, v \in W$, which have the form $u = g(z)$ and $v = h(z)$ for certain polynomials $g, h \in F[x]$ that are either zero or have degree less than d . Note $uv = g(z)h(z) = (gh)(z)$, where the product polynomial $gh \in F[x]$ may have degree as high as $2d - 2$. We will argue that for all $p \in F[x]$ of any degree, $p(z) \in W$; it will follow that $uv = (gh)(z)$ is indeed in W . First we use induction on n to show that $z^n \in W$ for all $n \geq 0$. This is true by definition of W for $0 \leq n < d$. Assume $n \geq d$ and we already know that $z^0, z^1, \dots, z^{n-1} \in W$. Let the minimal polynomial m of z over F be $m = a_0 + a_1x + \dots + a_{d-1}x^{d-1} + x^d$, where all $a_i \in F$. Multiplying by x^{n-d} and evaluating at $x = z$, we find that

$$z^n = -a_0z^{n-d} - a_1z^{n-d+1} - \dots - a_{d-1}z^{n-1}.$$

By induction, each power of z on the right side is in W . So z^n , which is an F -linear combination of these powers, is also in W . It now follows that $p(z)$, which is an F -linear combination of powers of z , is also in W .

Finally, is W closed under multiplicative inverses? Fix a nonzero element $u \in W$, which has the form $u = g(z)$ for some nonzero $g \in F[x]$ of degree less than d . As m is irreducible and $\deg(g) < \deg(m)$, we must have $\gcd(g, m) = 1$, and so $1 = gr + ms$ for some $r, s \in F[x]$. Evaluating both sides at z gives $1_K = g(z)r(z) + m(z)s(z) = ur(z)$. Now, as seen in the previous paragraph, $r(z) \in W$; and $ur(z) = 1$ shows that $u^{-1} = r(z) \in W$. This completes the proof that W is a subfield, and hence that $W = F(z)$.

We have now proved that *if z is algebraic over F , then the degree $[F(z) : F]$ is finite and equals the degree of z 's minimal polynomial over F .* Conversely, if $d = [F(z) : F]$ is finite, then z is algebraic over F (and so z 's minimal polynomial exists and has degree d). To see this, assume $[F(z) : F] = d < \infty$, and consider the list of $d + 1$ vectors $(1, z, z^2, \dots, z^d)$ in the F -vector space $F(z)$. Since $F(z)$ is d -dimensional, this list must be linearly dependent over F (see §1.8). So there exist scalars $c_0, c_1, \dots, c_d \in F$, not all zero, with $\sum_{i=0}^d c_i z^i = 0_K$. Then the polynomial $\sum_{i=0}^d c_i x^i \in F[x]$ is nonzero and has z as a root, so z is algebraic over F as claimed.

Here is an example that will be used in the discussion of SQC below. Suppose F is a subfield of K , z is in K but not in F , and $y = z^2$ is in F . Then $[F(z) : F] = 2$ and $(1, z)$ is an ordered F -basis of the F -vector space $F(z)$. Indeed, since $x^2 - z^2 = x^2 - y$ is a nonzero polynomial in $F[x]$ having z as a root, we see that z is algebraic over F . The polynomial $x^2 - y = (x - z)(x + z)$ does not factor into linear factors in $F[x]$, since z is not in F . It follows that $x^2 - y$ is irreducible in $F[x]$; it is also monic and has z as a root, so this is the minimal polynomial of z over F . The stated conclusions now follow from the results proved above, since the minimal polynomial has degree 2. For instance, taking $F = \mathbb{Q}$ and $z = \sqrt{2}$, we see that $\mathbb{Q}(\sqrt{2})$ is a two-dimensional \mathbb{Q} -vector space with ordered basis $(1, \sqrt{2})$. So, every element of the field $\mathbb{Q}(\sqrt{2})$ has a unique representation in the form $a + b\sqrt{2}$ with $a, b \in \mathbb{Q}$.

12.8 Proof that $\text{AC} = \text{SQC}$

We are now ready to prove that $\text{AC} \subseteq \text{SQC}$, i.e., that every arithmetically constructible complex number has a square root tower. Let $z \in \text{AC}$ have an arithmetic construction sequence z_0, z_1, \dots, z_k . We use this sequence to build a square root tower for z , as follows. Start with $\mathbb{Q} = K_0$, which contains $z_0 = 0$ and $z_1 = 1$. Suppose inductively that, for some $i \leq k$, we already have a chain of fields

$$\mathbb{Q} = K_0 \subseteq K_1 \subseteq K_2 \subseteq \cdots \subseteq K_s \subseteq \mathbb{C}$$

such that each $[K_j : K_{j-1}] = 2$ and all the elements z_0, z_1, \dots, z_{i-1} lie in K_s . Consider the possible cases for z_i . If z_i is the sum or product of two preceding elements in the sequence, or if z_i is the additive inverse or multiplicative inverse of a preceding element, then z_i lies in K_s because K_s is a subfield of \mathbb{C} . The only other possibility is that $z_i^2 = z_r$ for some $r < i$. If z_i happens to lie in K_s already, we can keep the current chain of fields. Otherwise, the example at the end of the last section applies to show that $K_{s+1} = K_s(z_i)$ is a field extension of K_s that contains z_i and satisfies $[K_{s+1} : K_s] = 2$. Thus, the inductive construction can continue. Proceeding in this way for $i = 2, 3, \dots, k$, we eventually obtain a square root tower for z , so that $z \in \text{SQC}$.

To prove $\text{SQC} \subseteq \text{AC}$, we need one final lemma about field extensions. Suppose $F \subseteq K$ are subfields of \mathbb{C} such that $[K : F] = 2$. We assert that there is a $w \in K \setminus F$ such that $w^2 \in F$, $K = F(w)$, and $(1, w)$ is an ordered basis for the F -vector space K . Start by picking any $v \in K \setminus F$. We have a chain of fields $F \subseteq F(v) \subseteq K$, so the degree formula gives $2 = [K : F] = [K : F(v)][F(v) : F]$. Since $v \notin F$, the field $F(v)$ properly contains F , which forces $[F(v) : F] = 2$, $[K : F(v)] = 1$, and hence $K = F(v)$. The minimal polynomial of v exists and has degree $2 = [F(v) : F]$; let this minimal polynomial be $x^2 + bx + c$ with $b, c \in F$. We know $v^2 + bv + c = 0_K$; completing the square gives $(v + b/2)^2 = b^2/4 - c \in F$. Now $w = v + b/2 \in K$, $w^2 \in F$, and $w \notin F$, since otherwise $v = w - b/2$ would lie in F . The

same argument used for v shows that $K = F(w)$ and $[F(w) : F] = 2$. By previous results, $(1, w)$ is an F -basis of K , so every element of K can be written uniquely in the form $a + dw$ for some $a, d \in F$.

Now suppose $z \in \text{SQC}$, and let

$$\mathbb{Q} = K_0 \subseteq K_1 \subseteq K_2 \subseteq \cdots \subseteq K_t \subseteq \mathbb{C}$$

be a square root tower for z . We will show by induction on j that every $K_j \subseteq \text{AC}$ for $0 \leq j \leq t$, which implies $z \in \text{AC}$ since $z \in K_t$. First consider $K_0 = \mathbb{Q}$. Since $0 \in \text{AC}$, $1 \in \text{AC}$, and AC is closed under addition, a quick induction shows that every positive integer n is in AC . So all integers are in AC by A2, hence all rational numbers are in AC by A3 and A4. For the induction step, assume $0 \leq j < t$ and $K_j \subseteq \text{AC}$. Use the lemma to get $w \in K_{j+1}$ such that $w^2 \in K_j$ and every element of K_{j+1} has the form $a + dw$ for some $a, d \in K_j$. Now $w^2 \in K_j \subseteq \text{AC}$, so A5 shows that $w \in \text{AC}$. Then the induction hypothesis, A1, and A3 show that $a + dw \in \text{AC}$ for all $a, d \in K_j$. So $K_{j+1} \subseteq \text{AC}$, completing the induction.

12.9 Impossibility of Geometric Construction Problems

In this section, we apply (12.2) to prove that the problems of duplicating a cube, squaring a circle, inscribing a regular heptagon in a circle, and trisecting a 60° angle are unsolvable with ruler and compass. To use (12.2), we will compute $[\mathbb{Q}(z) : \mathbb{Q}]$ by finding the degree of the minimal polynomial of z over \mathbb{Q} . In turn, to check that a given polynomial $m \in \mathbb{Q}[x]$ is the minimal polynomial of a given $z \in \mathbb{C}$, it suffices to verify that m is monic, $m(z) = 0$, and m is irreducible in $\mathbb{Q}[x]$. The following results (proved in Chapter 3) can help establish the irreducibility of m .

First, given $m \in \mathbb{Q}[x]$, there is a nonzero $c \in \mathbb{Q}$ with $cm \in \mathbb{Z}[x]$, and m is irreducible in $\mathbb{Q}[x]$ iff cm is irreducible in $\mathbb{Q}[x]$. So we need only test irreducibility of polynomials with integer coefficients. If the new polynomial is $cm = a_n x^n + \cdots + a_1 x + a_0$ (each $a_i \in \mathbb{Z}$), the rational root theorem states that any rational root of cm (or of m) must have the form p/q , where p divides a_0 and q divides a_n in \mathbb{Z} . Thus, we can find all rational roots of our polynomial by checking finitely many possibilities. Next, a polynomial of degree 2 or 3 in $\mathbb{Q}[x]$ with no rational roots is automatically irreducible in $\mathbb{Q}[x]$.

Another criterion for irreducibility in $\mathbb{Q}[x]$ involves reduction of the coefficients modulo a prime p . Suppose $cm = a_n x^n + \cdots + a_1 x + a_0 \in \mathbb{Z}[x]$ where a_n is not divisible by p . Consider the polynomial $q = (a_n \bmod p)x^n + \cdots + (a_1 \bmod p)x + (a_0 \bmod p) \in \mathbb{Z}_p[x]$ obtained by reducing each integer coefficient modulo p . If the polynomial q is irreducible in $\mathbb{Z}_p[x]$ for some prime p , then the original polynomial cm is irreducible in $\mathbb{Q}[x]$. Now, irreducibility in the polynomial ring $\mathbb{Z}_p[x]$ can be checked algorithmically since there are only finitely many possible factors for q . Unfortunately, this sufficient condition for irreducibility is not necessary: one can find irreducible polynomials in $\mathbb{Q}[x]$ whose reductions modulo *every* prime are reducible. Luckily, there is an algorithm (called *Kronecker's method*) for testing in finitely many steps whether an arbitrary polynomial $cm \in \mathbb{Z}[x]$ is irreducible in $\mathbb{Q}[x]$; see §3.18.

Now we are ready to tackle the construction problems mentioned at the beginning of this section. First consider duplication of the cube. Given a cube with side length 1, the cube with twice the volume has side length $\sqrt[3]{2}$. One readily checks that if there exist $P, Q \in \text{GC}$ such that the line segment \overline{PQ} has length $\sqrt[3]{2}$, then the real number $\sqrt[3]{2}$ would lie in GC . But the minimal polynomial of $\sqrt[3]{2}$ is $x^3 - 2 \in \mathbb{Q}[x]$, which is a degree 3 polynomial

with no rational roots (since $\pm 1, \pm 2$ are not roots), hence is irreducible in $\mathbb{Q}[x]$. Now $[\mathbb{Q}(\sqrt[3]{2}) : \mathbb{Q}] = \deg(x^3 - 2) = 3$ is not a power of 2. By (12.2), $\sqrt[3]{2}$ is not in GC.

Next consider the problem of trisecting the angle $\theta = 60^\circ$. We can assume that the angle in question is formed by the x -axis and the line $y = x\sqrt{3}$ passing through $(\cos 60^\circ, \sin 60^\circ)$. If this angle could be trisected, then $(\cos 20^\circ, \sin 20^\circ) \in \text{GC}$, so (by CR5) the real number $\alpha = \cos 20^\circ$ would be in GC. We can compute the minimal polynomial of α in $\mathbb{Q}[x]$ as follows. First, repeated use of the trigonometric addition formulas leads to the identity $\cos(3t) = 4(\cos t)^3 - 3\cos t$. Letting $t = 20^\circ$, we see that $1/2 = 4\alpha^3 - 3\alpha$. It follows that α is a root of the monic polynomial $x^3 - (3/4)x - 1/8 \in \mathbb{Q}[x]$, as well as the polynomial $8x^3 - 6x - 1$ with integer coefficients. The possible rational roots of these polynomials are $\pm 1, \pm 1/2, \pm 1/4$, and $\pm 1/8$, none of which work. Since these polynomials have degree 3 with no roots in \mathbb{Q} , they are irreducible in $\mathbb{Q}[x]$. So $x^3 - (3/4)x - 1/8$ is the minimal polynomial of α over \mathbb{Q} , $[\mathbb{Q}(\alpha) : \mathbb{Q}] = 3$, and we see $\cos 20^\circ \notin \text{GC}$ by (12.2).

Now, consider the problem of inscribing a regular heptagon (seven-sided polygon) in the unit circle. If one vertex of the heptagon is $(1, 0)$, the next vertex will be $(\cos(2\pi/7), \sin(2\pi/7)) = e^{2\pi i/7}$. So we must show that $\omega = e^{2\pi i/7} \notin \text{GC}$. What is the minimal polynomial of ω in $\mathbb{Q}[x]$? On one hand, ω is certainly a root of the polynomial $x^7 - 1$. But this polynomial is reducible, with factorization

$$x^7 - 1 = (x - 1)(x^6 + x^5 + x^4 + x^3 + x^2 + x + 1).$$

Since ω is not a root of $x - 1$, it must be a root of $m = x^6 + x^5 + \dots + x + 1 \in \mathbb{Z}[x]$. One may check that the reduction of m in $\mathbb{Z}_3[x]$ is irreducible, and hence m is irreducible in $\mathbb{Q}[x]$. (More generally, for any prime p , it can be shown that $1 + x + x^2 + \dots + x^{p-1}$ is irreducible in $\mathbb{Q}[x]$; see Exercise 66 in Chapter 3.) Thus, m is the minimal polynomial of ω . Since m has degree six, $[\mathbb{Q}(\omega) : \mathbb{Q}] = 6$. Six is not a power of 2, so ω is not geometrically constructible.

Finally, consider the notorious problem of squaring the circle. A circle with radius 1 has area π , so the square with the same area would have side length $\sqrt{\pi}$. If $\sqrt{\pi}$ were in GC, we would also have $\pi \in \text{GC}$ by CR8. But a hard theorem of Lindemann states that π is not algebraic over \mathbb{Q} . By the contrapositive of a result from §12.7, we see that $[\mathbb{Q}(\pi) : \mathbb{Q}] = \infty$. Since this degree is not a (finite) power of 2, (12.2) applies to show that $\sqrt{\pi}$ and π are not in GC. We will not prove Lindemann's theorem here; the interested reader may find a proof in [30, Vol. 1, Chapter 4] or [22, §1.7].

12.10 Constructibility of the 17-Gon

In this section we show that a regular 17-sided polygon *can* be constructed using a ruler and compass. It suffices to show that $\omega = e^{2\pi i/17} \in \text{GC}$, since once the line segment from 1 to ω is available, the other segments can be readily constructed. Let $m = x^{16} + x^{15} + \dots + x + 1 \in \mathbb{Q}[x]$. By Kronecker's algorithm, or by reduction mod 3, or by Exercise 66 in Chapter 3, one can verify that m is irreducible in $\mathbb{Q}[x]$. Moreover, $x^{17} - 1 = (x - 1)m$, from which it follows that m is the minimal polynomial of ω over \mathbb{Q} . Thus, $[\mathbb{Q}(\omega) : \mathbb{Q}] = \deg(m) = 16 = 2^4$.

Unfortunately, this information alone is not enough to conclude that $\omega \in \text{GC}$, since the converse of (12.1) is not always true. So we adopt a different approach. Let $\theta = 2\pi/17$. We will develop an explicit formula for $\cos \theta$ that will show $\cos \theta \in \text{AC}$. Since $\text{AC} = \text{GC}$, one can construct $\cos \theta$ geometrically, and then one can construct $e^{2\pi i/17}$ by drawing the altitude to the real axis through $\cos \theta$.

We will make repeated use of the following observation. Suppose r_1 and r_2 are complex

numbers with $r_1 + r_2 = b$ and $r_1 r_2 = c$. Then $(x - r_1)(x - r_2) = x^2 - bx + c$, hence (by the quadratic formula) $r_1, r_2 = (b \pm \sqrt{b^2 - 4c})/2$.

Since $\omega^j = e^{2\pi ij/17} = \cos(j\theta) + i \sin(j\theta)$ and $\omega^{-j} = \cos(-j\theta) + i \sin(-j\theta) = \cos(j\theta) - i \sin(j\theta)$, we have $\omega^j + \omega^{-j} = 2 \cos(j\theta)$ for all $j \geq 1$. Consider the real numbers

$$x_1 = 2(\cos \theta + \cos(2\theta) + \cos(4\theta) + \cos(8\theta)); \quad x_2 = 2(\cos(3\theta) + \cos(5\theta) + \cos(6\theta) + \cos(7\theta)).$$

The previous remark shows that the sum $x_1 + x_2$ is

$$\omega + \omega^{-1} + \omega^2 + \omega^{-2} + \omega^4 + \omega^{-4} + \omega^8 + \omega^{-8} + \omega^3 + \omega^{-3} + \omega^5 + \omega^{-5} + \omega^6 + \omega^{-6} + \omega^7 + \omega^{-7}.$$

Since $\omega^{17} = 1$, we have $\omega^{-j} = \omega^{17-j}$, so we can also write the sum as

$$x_1 + x_2 = \sum_{i=1}^{16} \omega^i = m(\omega) - 1 = -1.$$

Next, recall the trigonometric identities $\cos(u+v) + \cos(u-v) = 2 \cos u \cos v$ and $\cos(-u) = \cos u$. Also, for $k > 8$, $\cos(k\theta) = \cos(-k\theta) = \cos((17-k)\theta)$ since $17\theta = 2\pi$. Using these identities along with the distributive law, we find that

$$\begin{aligned} x_1 x_2 &= 2(\cos(4\theta) + \cos(2\theta) + \cos(6\theta) + \cos(4\theta) + \cos(7\theta) + \cos(5\theta) + \cos(8\theta) + \cos(6\theta) \\ &\quad + \cos(5\theta) + \cos(1\theta) + \cos(7\theta) + \cos(3\theta) + \cos(8\theta) + \cos(4\theta) + \cos(8\theta) + \cos(5\theta) \\ &\quad + \cos(7\theta) + \cos(1\theta) + \cos(8\theta) + \cos(1\theta) + \cos(7\theta) + \cos(2\theta) + \cos(6\theta) + \cos(3\theta) \\ &\quad + \cos(6\theta) + \cos(5\theta) + \cos(4\theta) + \cos(3\theta) + \cos(3\theta) + \cos(2\theta) + \cos(2\theta) + \cos(1\theta)). \end{aligned}$$

Gathering terms, we discover that $x_1 x_2 = 4(x_1 + x_2) = -4$. By our initial observation, $x_1, x_2 = (-1 \pm \sqrt{17})/2$. Now, $x_1 \approx 1.56155 > 0$ and $x_2 \approx -2.56155 < 0$, so we conclude:

$$x_1 = \frac{-1 + \sqrt{17}}{2}, \quad x_2 = \frac{-1 - \sqrt{17}}{2}.$$

The next step is to consider the four real numbers

$$\begin{aligned} y_1 &= 2(\cos \theta + \cos(4\theta)), \quad y_2 = 2(\cos(2\theta) + \cos(8\theta)), \\ y_3 &= 2(\cos(3\theta) + \cos(5\theta)), \quad y_4 = 2(\cos(6\theta) + \cos(7\theta)). \end{aligned}$$

We see that $y_1 + y_2 = x_1$ and, by the same trigonometric identities used above,

$$y_1 y_2 = 2(\cos(3\theta) + \cos(1\theta) + \cos(8\theta) + \cos(7\theta) + \cos(6\theta) + \cos(2\theta) + \cos(5\theta) + \cos(4\theta)).$$

So $y_1 y_2 = x_1 + x_2 = -1$. The initial observation applies again to give

$$y_1 = \frac{x_1 + \sqrt{x_1^2 + 4}}{2}, \quad y_2 = \frac{x_1 - \sqrt{x_1^2 + 4}}{2},$$

where the signs were found by comparison to decimal approximations of y_1 and y_2 . Similarly, $y_3 + y_4 = x_2$ and

$$y_3 y_4 = 2(\cos(8\theta) + \cos(3\theta) + \cos(7\theta) + \cos(4\theta) + \cos(6\theta) + \cos(1\theta) + \cos(5\theta) + \cos(2\theta)).$$

So $y_3 y_4 = x_1 + x_2 = -1$, and hence

$$y_3 = \frac{x_2 + \sqrt{x_2^2 + 4}}{2}, \quad y_4 = \frac{x_2 - \sqrt{x_2^2 + 4}}{2}.$$

Finally, consider $z_1 = 2 \cos \theta$ and $z_2 = 2 \cos(4\theta)$. We have $z_1 + z_2 = y_1$ and $z_1 z_2 = 2 \cos(5\theta) + 2 \cos(3\theta) = y_3$. One final application of the initial observation gives

$$z_1 = \frac{y_1 + \sqrt{y_1^2 - 4y_3}}{2},$$

where the sign may be determined by noting that $z_1 > z_2 > 0$. We see by inspection of the formulas that $x_1, x_2 \in \text{AC}$, hence each $y_i \in \text{AC}$, hence $z_1 \in \text{AC}$, hence finally $\cos(2\pi/17) = z_1/2 \in \text{AC} = \text{GC}$. More explicitly, if we substitute formulas for y_1 , etc., into the formula for z_1 and simplify, we find that $\cos(2\pi/17)$ equals

$$\frac{-1 + \sqrt{17} + \sqrt{34 - 2\sqrt{17}} + \sqrt{\left(-1 + \sqrt{17} + \sqrt{34 - 2\sqrt{17}}\right)^2 + 16 \left(1 + \sqrt{17} - \sqrt{34 + 2\sqrt{17}}\right)}}{16},$$

which is evidently in AC!

12.11 Overview of Solvability by Radicals

The recursive definition of AC allows us to construct new numbers by using addition, subtraction, multiplication, division, and the extraction of square roots. If we expand this definition to allow the extraction of n 'th roots for any $n \geq 2$, we can obtain even more numbers. This leads to the following recursive definition of RC, the set of *radically constructible* complex numbers.

- R0. $0 \in \text{RC}$ and $1 \in \text{RC}$.
- R1. If $a, b \in \text{RC}$ then $a + b \in \text{RC}$.
- R2. If $a \in \text{RC}$ then $-a \in \text{RC}$.
- R3. If $a, b \in \text{RC}$ then $a \cdot b \in \text{RC}$.
- R4. If $a \in \text{RC}$ is nonzero, then $a^{-1} \in \text{RC}$.
- R5. If $a \in \text{RC}$, $b \in \mathbb{C}$, $n \geq 2$ is an integer, and $b^n = a$, then $b \in \text{RC}$.
- R6. The only numbers in RC are those that can be obtained by applying rules R0 through R5 a finite number of times.

Equivalently, $z \in \text{RC}$ iff there is a finite sequence

$$z_0 = 0, z_1 = 1, z_2, \dots, z_k = z \quad (k \geq 0)$$

such that, for $2 \leq i \leq k$, either z_i is the sum or product of two preceding numbers in the sequence, or z_i is the additive inverse, multiplicative inverse, or n 'th root of some preceding number in the sequence. Such a sequence is called a *radical construction sequence* for z .

For example, consider a quadratic polynomial $ax^2 + bx + c$ with $a, b, c \in \text{RC}$ and $a \neq 0$. The roots of this polynomial are $(-b \pm \sqrt{b^2 - 4ac})/(2a)$, which are also in RC. Next, consider a “reduced” cubic polynomial $y^3 + qy + r$ where $q, r \in \text{RC}$. The *cubic formula* states that the roots of this polynomial have the form $u - q/(3u)$ where u is any complex number satisfying $u^3 = (-r + \sqrt{r^2 + 4q^3}/27)/2$. It follows that these roots are also in RC. Similarly, the *quartic formula* can be used to show that the roots of a fourth-degree polynomial with coefficients in RC are also in RC. However, there exist fifth-degree polynomials in $\mathbb{Z}[x]$ whose roots do not lie in RC. This is the famous theorem that quintic (degree 5) polynomials are not always “solvable by radicals.”

We only have space here to give an outline of the proof strategy. The first step is to give a field-theoretic rephrasing of the definition of RC, which is analogous to the result AC = SQC proved above. Specifically, for a prime p , let us say that a field extension $F \subseteq K$ is *pure* of type p if there exists $w \in K$ with $K = F(w)$ and $w^p \in F$; informally, we obtain K from F by “adjoining a p ’th root of an element of F .” One can show that $z \in \text{RC}$ iff there is a chain of fields

$$\mathbb{Q} = K_0 \subseteq K_1 \subseteq K_2 \subseteq \cdots \subseteq K_t \subseteq \mathbb{C}$$

where $z \in K_t$ and each K_j is a pure extension of K_{j-1} .

The second step is to use *Galois theory* to convert a chain of fields of the type just described to a chain of finite groups

$$G = G_0 \supseteq G_1 \supseteq G_2 \supseteq \cdots \supseteq G_t = \{e\}$$

in which G_i is a normal subgroup of G_{i-1} for each $i \geq 1$, and G_{i-1}/G_i is a cyclic group of size p . Not all finite groups G possess chains of subgroups satisfying these conditions; those groups that do are called *solvable*. Finally, to each polynomial in $\mathbb{Q}[x]$ one can associate a field extension of \mathbb{Q} (the “splitting field” of the polynomial) and a finite group (the “Galois group” of the polynomial). One can construct degree 5 polynomials whose associated Galois groups are not solvable, and this can be combined with the preceding results to show that the roots of these polynomials cannot all lie in RC. For more information, consult the texts on Galois theory listed in the bibliography.

12.12 Summary

1. *Geometrically Constructible Numbers.* GC is the set of all points P in the complex plane for which there is a sequence $P_0 = 0, P_1 = 1, P_2, \dots, P_k = P$ such that each P_j ($j \geq 2$) is an intersection point of two lines and/or circles determined by preceding points in the sequence.
2. *Arithmetically Constructible Numbers.* AC is the set of all complex numbers z for which there is a sequence $z_0 = 0, z_1 = 1, z_2, \dots, z_k = z$ such that each z_j ($j \geq 2$) is the sum, additive inverse, product, multiplicative inverse, or square root of preceding number(s) in the sequence.
3. *Field Extensions.* Given a field K and subfield F , the degree $[K : F]$ is the dimension of K viewed as a vector space over F . For $z \in K$, the subfield $F(z)$ is the intersection of all subfields of K containing F and z . Given a chain of fields $F_0 \subseteq F_1 \subseteq F_2 \subseteq \cdots \subseteq F_n$, we have the degree formula $[F_n : F_0] = \prod_{i=1}^n [F_i : F_{i-1}]$.
4. *Square Root Towers.* SQC is the set of all complex numbers z for which there is a chain of fields $\mathbb{Q} = K_0 \subseteq K_1 \subseteq \cdots \subseteq K_t \subseteq \mathbb{C}$ with $z \in K_t$ and $[K_i : K_{i-1}] = 2$ for all $i \geq 1$. We have $K_i = K_{i-1}(w)$ for some $w \in K_i$ such that $w^2 \in K_{i-1}$, and $K_i = \{a + bw : a, b \in K_{i-1}\}$. These facts are used to show $\text{SQC} = \text{AC}$.
5. *Criteria for Constructibility.* $\text{GC} = \text{AC} = \text{SQC}$. For a point $z \in \mathbb{C}$ to be constructible by ruler and compass, it is necessary (but not always sufficient) that $[\mathbb{Q}(z) : \mathbb{Q}]$ be a finite power of 2. In particular, $\sqrt[3]{2}$, π , $\cos 20^\circ$, and $e^{2\pi i/7}$ are not constructible, but $e^{2\pi i/17}$ is constructible.
6. *Algebraic Elements and Minimal Polynomials.* Given a field K , a subfield F , and

$z \in K$, z is algebraic over F iff z is the root of some nonzero polynomial in $F[x]$; this holds iff $[F(z) : F]$ is finite. In this case, there exists a unique monic irreducible polynomial $m \in F[x]$ having z as a root, called the minimal polynomial of z over F . Letting $d = \deg(m)$, $(1, z, z^2, \dots, z^{d-1})$ is an ordered F -basis of $F(z)$, so that $[F(z) : F] = d$, the degree of z 's minimal polynomial.

7. *Proof Idea for $AC \subseteq GC$.* AC is contained in GC because it is possible to simulate each arithmetical operation $+$, $-$, \times , \div , $\sqrt{}$ by geometrical constructions. For example, if we draw the circle with center $(x - 1)/2$ passing through $x \in \mathbb{R}^+$, this circle hits the positive imaginary axis at the point $i\sqrt{x}$. Similar triangles can be used to implement real multiplication and division. Altitudes and parallels can be used to transfer distances and thereby perform real addition and subtraction. Finally, complex arithmetic can be reduced to real arithmetic on the real and imaginary parts, and complex square roots of $re^{i\theta}$ can be found by taking the square root of the real modulus r and bisecting the angle θ .
 8. *Proof Idea for $GC \subseteq AC$.* GC is contained in AC because the coordinates of the intersection points of two lines and/or circles can be determined from the coordinates of points determining those lines and circles by arithmetic operations and square root extractions.
-

12.13 Exercises

1. Give explicit arithmetic construction sequences for each of these complex numbers: (a) 5; (b) $-2 + i$; (c) $\sqrt{7}$; (d) $2e^{5\pi i/6}$.
2. Give explicit geometric construction sequences for each of these complex numbers: (a) 5; (b) $-2 + i$; (c) $\sqrt{7}$; (d) $2e^{5\pi i/6}$.
3. Give explicit square root towers for each of these complex numbers: (a) 5; (b) $-2 + i$; (c) $\sqrt{7}$; (d) $2e^{5\pi i/6}$.
4. (a) Find all $z \in \mathbb{C}$ that have geometric construction sequences of length at most 3. (b) Find all $z \in \mathbb{C}$ that have arithmetic construction sequences of length at most 4.
5. (a) Find all $z \in \mathbb{C}$ that have geometric construction sequences of length at most 4. (b) Find all $z \in \mathbb{C}$ that have arithmetic construction sequences of length at most 5. (You may want to use a computer for this problem.)
6. (a) Use the definition of AC to prove in detail that $\mathbb{Q} \subseteq AC$. (b) Use the definition of AC to prove that AC is a subfield of \mathbb{C} .
7. Let n be a positive integer. (a) Prove that $0, 1, 2, \dots, n$ is a geometric construction sequence for n . (b) Prove: if $P_0 = 0, P_1 = 1, P_2, \dots, P_k = n$ is any geometric construction sequence for n , then $P_0, \dots, P_k, 2n$ is a geometric construction sequence for $2n$. (c) Prove: if $P_0 = 0, P_1 = 1, P_2 = -1, \dots, P_k = n$ is a geometric construction sequence for n , then $P_0, \dots, P_k, 2n + 1$ is a geometric construction sequence for $2n + 1$. (d) Find a geometric construction sequence for 99 of length 9. (e) Prove: if n is k bits long when written in base 2, then n has a geometric construction sequence of length $k + 2$.
8. Prove: for all positive integers n that are k bits long when written in base 2, n

- has an arithmetic construction sequence of length at most $2k$. (Use induction on k .)
9. (a) Let $a, b \in \text{AC}$ have arithmetic construction sequences $x_0, \dots, x_k = a$ and $y_0, \dots, y_m = b$. Prove that $x_0, \dots, x_k, y_0, \dots, y_m, a + b$ is an arithmetic construction sequence for $a + b$. (b) Prove that the recursive definition of AC (using rules A0 through A6) is logically equivalent to the iterative definition of AC (using arithmetic construction sequences).
 10. Let $\theta \in \mathbb{R}$. (a) Give a geometric proof that $\cos \theta \in \text{GC}$ iff $\sin \theta \in \text{GC}$ iff $e^{i\theta} \in \text{GC}$. (b) Give an arithmetic proof that $\cos \theta \in \text{AC}$ iff $\sin \theta \in \text{AC}$ iff $e^{i\theta} \in \text{AC}$. (c) True or false: for all θ , if $e^{i\theta} \in \text{AC}$ then $\theta \in \text{AC}$.
 11. Let K be a field with subfield F . Prove that K is an F -vector space by writing out the field axioms for K and comparing them to the axioms for a vector space over F .
 12. (a) Let K be a field with subfield F . Prove that $F = K$ iff $[K : F] = 1$. (b) Give an example of a chain of fields $F \subseteq K \subseteq E$ where $[E : F] = [K : F]$, yet $K \neq E$.
 13. Let $F = \{a + b\sqrt{3} + c\sqrt{5} + d\sqrt{15} : a, b, c, d \in \mathbb{Q}\}$. (a) Prove F is a subfield of \mathbb{C} . (b) Compute $[F : \mathbb{Q}]$, $[F : \mathbb{Q}(\sqrt{3})]$, and $[F : \mathbb{Q}(\sqrt{5})]$, and find ordered bases for each of these vector spaces. (c) Find five different subfields of F , and prove carefully that they are all different.
 14. Suppose $F \subseteq E \subseteq K$ is a chain of fields with $[K : F] = p$, a prime integer. Prove that $E = F$ or $E = K$.
 15. (a) Given a chain of fields $F \subseteq K \subseteq E$, prove that $[E : F]$ is finite iff $[E : K]$ is finite and $[K : F]$ is finite. (b) State and prove a similar result for a chain of more than three fields.
 16. *Field Extension Generated by a Subset.* Let K be a field with subfield F , and let S be an arbitrary subset of K . Define $F(S)$ to be the intersection of all subfields of K containing $F \cup S$. (a) Prove $F(S)$ is a subfield of K , $F \subseteq F(S)$, and $S \subseteq F(S)$. (b) Prove: for any subfield M of K , $F(S) \subseteq M$ iff $F \subseteq M$ and $S \subseteq M$. (c) Prove: for all subsets S, T of K , $F(S)(T) = F(T)(S) = F(S \cup T)$. (d) Prove: for all subsets S, T of K , $F(S) = F(T)$ iff $S \subseteq F(T)$ and $T \subseteq F(S)$. (e) Give a specific example of F and K and two disjoint nonempty subsets S and T of $K \sim F$ with $|S| \neq |T|$ and yet $F(S) = F(T)$.
 17. Let K be a field. (a) Prove that the set $\{n.1_K : n \in \mathbb{Z}\}$ of additive multiples of 1 in K is a subring of K that is isomorphic either to \mathbb{Z} or to \mathbb{Z}_p for some prime p . (b) Prove that the intersection of all subfields of K is a field that is isomorphic either to \mathbb{Q} or to \mathbb{Z}_p for some prime p .
 18. Let K be a field with subfield F , and fix $z \in K$. (a) Prove $F(z) = \{f(z)/g(z) : f, g \in F[x] \text{ and } g(z) \neq 0\}$. [Hint: Call the right-hand side E . To prove $F(z) \subseteq E$, show E is a subfield of K containing F and z . To prove $E \subseteq F(z)$, show that every subfield M of K containing F and z must also contain E .] (b) Let $R = \{f(z) : f \in F[x]\}$. Prove that R is a subring of K that contains F and z . Prove that if z is algebraic over F , then $R = F(z)$. (c) Use the theorem that π is not algebraic over \mathbb{Q} to prove that $\{f(\pi) : f \in \mathbb{Q}[x]\}$ is not a subfield of \mathbb{C} .
 19. Describe a specific square root tower for $e^{2\pi i/17}$.
 20. Suppose $x_0, y_0, x_1, y_1 \in \text{AC}$ satisfy $x_0 \neq x_1$. Let the unique line through (x_0, y_0) and (x_1, y_1) have slope m and y -intercept b for some $m, b \in \mathbb{R}$. Prove that $m, b \in \text{AC}$.

21. Consider a line $y = mx + d$ and a parabola $y = ax^2 + bx + c$ where $a, b, c, d, m \in \text{AC}$. Prove that all intersection points of this line and this parabola lie in AC.
22. In the proof of CR1, use theorems about triangles to confirm the assertions that E is the midpoint of \overline{AB} and \overline{CD} is perpendicular to \overline{AB} .
23. Use an actual ruler and compass to construct: (a) a regular hexagon inscribed in the unit circle with one vertex at $(1, 0)$; (b) a square inscribed in the unit circle with one vertex at $e^{\pi i/4}$.
24. Use an actual ruler and compass to construct: (a) the line $y = 3x + 1$; (b) the altitude to the line in (a) passing through $(2, 0)$; (c) the line parallel to the line in (a) passing through $(2, 0)$.
25. Use an actual ruler and compass to construct: (a) the real numbers $x = \sqrt{6}$ and $y = -3/2$; (b) $x + y$; (c) $x - y$; (d) xy ; (e) $1/x$.
26. Use an actual ruler and compass to construct: (a) the complex numbers $u = 2 + i$ and $v = 3 - 2i$; (b) $u + v$; (c) $u - v$; (d) uv ; (e) the two complex square roots of u ; (f) $1/v$.
27. (a) Prove geometrically: if P and Q are in GC and r is the length of the line segment \overline{PQ} , then the real number $r = (r, 0)$ is in GC. Conclude that the collapsing compass and ruler can simulate a real compass's ability to transfer a fixed distance from one part of a diagram to another. (b) Reprove (a) algebraically using the fact that $\text{GC} = \text{AC}$.
28. (a) Prove geometrically: for all $\alpha, \beta \in \mathbb{R}$, if $e^{i\alpha} \in \text{GC}$ and $e^{i\beta} \in \text{GC}$, then $e^{i(\beta-\alpha)} \in \text{GC}$. Conclude that when deciding the constructibility of regular n -gons in the unit circle, it suffices to consider n -gons with one vertex at $(1, 0)$. (b) Reprove (a) algebraically using the fact that $\text{GC} = \text{AC}$.
29. (a) In CR11, we proved that for all complex $u, v \in \text{GC}$, $u + v \in \text{GC}$ by looking at real and imaginary parts. Give an alternate geometric proof of this result by using the parallelogram law for finding the vector sum of u and v in \mathbb{R}^2 . (b) Similarly, prove geometrically that $u, v \in \text{GC}$ implies $uv \in \text{GC}$ by using the polar multiplication formula $(re^{i\alpha}) \cdot (se^{i\beta}) = (rs)e^{i(\alpha+\beta)}$.
30. Given two points P and Q , show how to use a ruler and compass to construct points R and S that trisect the line segment \overline{PQ} . Illustrate your construction with an actual ruler and compass.
31. Consider the three points $A = e^{\pi i/3}$, $B = 0$, and $C = 1$, so $\angle ABC$ is a 60° angle. As shown in the previous problem, we can find points D and E that trisect the line segment \overline{AC} . (a) Write D and E in the form $a + bi$ with $a, b \in \mathbb{R}$. (b) Use (a) to approximate the measure of the angle $\angle EBC$ in degrees. Has this construction trisected $\angle ABC$?
32. Starting with the points 0 and 1, prove that a compass alone can be used to construct all points in the set $\{a + be^{2\pi i/6} : a, b \in \mathbb{Z}\}$. (Remarkably, it can be shown that *every* point in GC can be constructed with just a compass, starting from the two points 0 and 1.)
33. Find (with proof) the minimal polynomials over \mathbb{Q} of each of these complex numbers: (a) $22/7$; (b) $2i$; (c) $\sqrt{11}$; (d) $\sqrt{3} + \sqrt{7}$; (e) $\sqrt[5]{2}$; (f) $\sqrt{2} + i$; (g) $e^{2\pi i/6}$; (h) $e^{2\pi i/12}$.
34. Let F be a field, and let g be an irreducible polynomial in the polynomial ring $F[x]$. Let $I = F[x]g$ be the principal ideal in $F[x]$ generated by g , and let K be the quotient ring $F[x]/I$. (See Chapter 1 for the relevant definitions.) (a) Prove

that K is a field, and $F_1 = \{c + I : c \in F\}$ is a subfield of K isomorphic to F . (b) Use the isomorphism between F and F_1 to convert g to a polynomial $g_1 \in F_1[x]$. Prove that $g_1(x + I) = 0 + I$, so that the coset $x + I$ is a root of g_1 in the field K . (c) Suppose F is a subfield of \mathbb{C} , and $g \in F[x]$ is the minimal polynomial over F of some $z \in \mathbb{C}$. Prove that the field K in part (a) is isomorphic to the subfield $F(z) \subseteq \mathbb{C}$. (Apply the fundamental homomorphism theorem for rings to an appropriate map.)

35. Construct a specific field K with 343 elements containing an element z with $z^3 = 2$ and containing a subfield isomorphic to \mathbb{Z}_7 . (Use the previous exercise.)
36. Consider the polynomial $x^4 + x^3 + x^2 + x + 1 = (x^5 - 1)/(x - 1)$. (a) Find the complex roots of this polynomial in terms of cosines, sines, and/or complex exponentials. (b) Divide the polynomial by x^2 and let $y = x + x^{-1}$. Convert the quartic polynomial in x to a quadratic polynomial in y , and then use the quadratic formula to solve for y and then for x algebraically. (c) Use (a) and (b) to prove that

$$\cos 72^\circ = \frac{-1 + \sqrt{5}}{4}, \quad \sin 72^\circ = \frac{1}{2} \sqrt{\frac{5 + \sqrt{5}}{2}},$$

and hence that $e^{2\pi i/5} \in \text{GC}$. (d) Use an actual ruler and compass to construct a regular pentagon inscribed in the unit circle.

37. Let z be the unique real root of the polynomial $x^3 + 2x + 1 \in \mathbb{Q}[x]$. (a) Find an ordered \mathbb{Q} -basis of $\mathbb{Q}(z)$ and compute $[\mathbb{Q}(z) : \mathbb{Q}]$. (b) For $-1 \leq n \leq 7$, express z^n as a \mathbb{Q} -linear combination of the basis in (a). (c) Compute $(z^2 - 3z + 2) \cdot (2z^2 + 5z - 1)$, expressing the answer in terms of the basis in (a). (d) Find $(z^2 - z + 1)^{-1}$, expressing the answer in terms of the basis in (a). (For (b) through (d), use ideas from the proof in §12.7.)
38. Let $z = e^{2\pi i/12}$. (a) Find an ordered \mathbb{Q} -basis of $\mathbb{Q}(z)$ and compute $[\mathbb{Q}(z) : \mathbb{Q}]$. (b) For all $n \in \mathbb{Z}$, express z^n as a \mathbb{Q} -linear combination of the basis in (a). (c) Compute $(z^3 + 2z - 5) \cdot (z^2 + 3z + 4)$, expressing the answer in terms of the basis in (a). (d) Find $(z^3 + 2z^2 + 3z + 4)^{-1}$, expressing the answer in terms of the basis in (a).
39. Suppose $z \in \mathbb{C}$ has a minimal polynomial over \mathbb{Q} of odd degree. Prove that $x^2 - 2$ is the minimal polynomial of $\sqrt{2}$ over the subfield $\mathbb{Q}(z)$.
40. (a) Let β be an irrational real number, and assume that for all integers $k \geq 1$ there exists a rational number $s_k = m_k/n_k$ (for some $m_k, n_k \in \mathbb{Z}$ with $n_k > 0$) satisfying $|\beta - s_k| < (kn_k^k)^{-1}$. Prove β is not algebraic over \mathbb{Q} , hence $\beta \notin \text{GC}$. (b) Deduce from (a) that $\beta = \sum_{n \geq 1} 10^{-n!}$ is not algebraic over \mathbb{Q} .
41. (a) Prove: for all $z, w \in \mathbb{C}$, if z and w are algebraic over \mathbb{Q} , then $z + w$ and zw are algebraic over \mathbb{Q} . [Hint: Consider the chain of fields $\mathbb{Q} \subseteq \mathbb{Q}(z) \subseteq \mathbb{Q}(z)(w)$ and look at degrees.] (b) Let K be the set of all $z \in \mathbb{C}$ that are algebraic over \mathbb{Q} . Prove that K is a subfield of \mathbb{C} containing \mathbb{Q} . (c) Is $[K : \mathbb{Q}]$ finite or infinite? Is K countable or uncountable?
42. (a) Suppose K is a field with subfield F such that $[K : F] = 2$ and $1_K + 1_K \neq 0_K$. Prove there exists $w \in K$ with $K = F(w)$ and $w^2 \in F$. (b) Let $K = \{0, 1, 2, 3\}$,

and define addition and multiplication in K by the following tables:

$+$	0	1	2	3	.	0	1	2	3
0	0	1	2	3	0	0	0	0	0
1	1	0	3	2	1	0	1	2	3
2	2	3	0	1	2	0	2	3	1
3	3	2	1	0	3	0	3	1	2

Prove that K is a field, $F = \{0, 1\}$ is a subfield of K , and $[K : F] = 2$. (c) Prove that the conclusion of (a) is false for the field K and subfield F in part (b).

43. For each n with $3 \leq n \leq 18$, decide (with proof) whether one can construct a regular n -gon with ruler and compass.
44. (a) Prove: for all $n \geq 3$, a regular n -gon is constructible with ruler and compass iff a regular $(2n)$ -gon is also constructible. (b) Prove that if a regular n -gon and a regular m -gon are constructible with ruler and compass and $k = \text{lcm}(m, n)$, then a regular k -gon is also constructible.
45. (a) Use the trigonometric addition formula for $\cos(\alpha + \beta)$ and other identities to derive the formula $\cos(3t) = 4(\cos t)^3 - 3\cos t$. (b) Similarly, prove that $\cos(u + v) + \cos(u - v) = 2\cos u \cos v$. (c) Reprove the formulas in (a) and (b) using $\cos t = (e^{it} + e^{-it})/2$ and facts about the exponential function.
46. Let $\theta \in (0, \pi/2)$ be the unique real number with $\cos \theta = 3/5$. (a) Find the minimal polynomial of $\theta/3$. (b) Explain how to construct $e^{i\theta}$ with a ruler and compass. (c) Prove or disprove: the angle between the x -axis and the ray through $e^{i\theta}$ can be trisected with ruler and compass.
47. Let $\theta \in (0, \pi/2)$ be the unique real number with $\cos \theta = 11/16$. (a) Find the minimal polynomial of $\theta/3$. (b) Prove or disprove: the angle between the x -axis and the ray through $e^{i\theta}$ can be trisected with ruler and compass.
48. In §12.10, verify that the signs in the formulas for y_1, y_2, y_3, y_4 are correct by calculating decimal approximations for these four quantities.
49. (a) Given a cubic polynomial $x^3 + bx^2 + cx + d$, show that the substitution $y = x + b/3$ leads to a new “reduced” cubic $y^3 + qy + r$ where the coefficient of y^2 is zero. (b) Find and justify a similar substitution that will change a monic degree n polynomial in x into a new degree n polynomial in y where the coefficient of y^{n-1} is zero.
50. *Cubic Formula.* Given $q, r \in \mathbb{C}$, let u be any nonzero complex solution of $u^3 = (-r + \sqrt{r^2 + 4q^3}/27)/2$. Verify by algebraic manipulations that $y = u - q/(3u)$ is a solution of $y^3 + qy + r = 0$. What happens in the case $u = 0$?
51. *Quartic Formula.* Find a formula for the roots of a “reduced” quartic polynomial $y^4 + py^2 + qy + r$ by hypothesizing a factorization of the form $(y^2 + ky + \ell)(y^2 - ky + m)$, obtaining equations for the unknown coefficients k, ℓ, m , eliminating ℓ and m , and finally obtaining a cubic equation in the variable k^2 . Solve this with the cubic formula and work back to find the roots of the original quartic. Conclude that if $p, q, r \in \text{RC}$, then all complex roots of the quartic polynomial are in RC .
52. Prove, as asserted in §12.11, that for all complex z , $z \in \text{RC}$ iff there is a chain of fields $\mathbb{Q} = K_0 \subseteq K_1 \subseteq K_2 \subseteq \cdots \subseteq K_t \subseteq \mathbb{C}$, where $z \in K_t$ and each K_j is a pure extension of K_{j-1} .

53. Prove that $\cos(2\pi/7) \in \text{RC}$ by finding an explicit algebraic formula for this number. [Hint: Use a substitution similar to the one in Exercise 36 to find the roots of $(x^7 - 1)/(x - 1)$.]
54. Find and justify an exact formula for $\cos 1^\circ$ that is built up from integers using only arithmetical operations, square roots, and cube roots. (By the trigonometric addition formulas, this shows there are “exact formulas” for $\cos k^\circ$ for any integer k . Curiously, grade school students are only made to memorize these formulas for a few selected choices of k .)
55. True or false? Explain each answer. (a) The minimal polynomial of $e^{2\pi i/3}$ over \mathbb{Q} is $x^3 - 1$. (b) GC is a subfield of \mathbb{C} . (c) AC = RC. (d) $e^{6\pi i/17} \in \text{GC}$. (e) For all integers $n \geq 1$, there is a chain of fields $\mathbb{Q} \subseteq K \subseteq \mathbb{C}$ with $[K : \mathbb{Q}] = n$. (f) For all $w, z \in \mathbb{C}$, $\mathbb{Q}(w) = \mathbb{Q}(z)$ iff $w = z$. (g) $e^{\pi i/17} \in \text{SQC}$. (h) For all fields K and subfields F with $[K : F] = 4$, there can exist at most one subfield E with $F \subsetneq E \subsetneq K$. (i) For every $k \geq 1$, there exists a sequence of distinct points P_0, P_1, \dots, P_k that is both an arithmetic construction sequence and a geometric construction sequence. (j) For any field K and any subfield F and any $z \in K$, if $[K : F]$ is finite then z is algebraic over F . (k) For every chain of fields $\mathbb{Q} \subseteq K \subseteq \mathbb{C}$ and all $z \in K$, if $[K : \mathbb{Q}] = 2$ then $z^2 \in \mathbb{Q}$. (l) For all real θ such that $0 < \theta < \pi/2$, if $e^{i\theta} \in \text{GC}$ then $e^{i\theta/3} \notin \text{GC}$.

This page intentionally left blank

13

Dual Spaces and Bilinear Forms

A recurring theme in mathematics is the relationship between *geometric spaces* and *structure-preserving functions* defined on these spaces. Galois theory for field extensions, elementary algebraic geometry, and the theory of dual vector spaces can all be viewed as instances of this common theme. Given a space V and a set R of functions on V , the basic idea is to study the *zero set* of a collection of functions in R , which is the set of points in V where all the functions in the collection are zero. Similarly, given a subset of V , one can consider the set of all functions in R that evaluate to zero at all points of the given subset. These two maps (one sending subsets of R to subsets of V , and the other sending subsets of V to subsets of R) often restrict to give a *one-to-one correspondence* between certain distinguished subsets of V and certain distinguished subsets of R . In many cases, key aspects of the geometric structure of V are mirrored by the algebraic structure of R through this correspondence. This interplay between V and R allows us to obtain structural information about the properties of both V and R .

This chapter studies one of the most basic and ubiquitous instances of this correspondence between spaces and functions, namely the theory of dual vector spaces. Our geometric space will be a finite-dimensional vector space V over a field F . The set R will consist of all *linear* functions mapping V into F . We will see that R is also a vector space over F (called the *dual space* for V and denoted V^*). The maps mentioned above will lead to an inclusion-reversing bijection between the subspaces of V and the subspaces of V^* . Later, we discuss the connection between these results and the theory of bilinear forms on V , real inner product spaces, and complex inner product spaces. After a brief interlude on Banach spaces, the chapter closes with a sketch of the “ideal-variety correspondence,” which is a cornerstone of algebraic geometry.

13.1 Vector Spaces of Linear Maps

Before commencing our study of dual spaces, we first give a general construction for manufacturing vector spaces. Let W be a vector space over a field F , and let S be any set. Let Z be the set of all functions $g : S \rightarrow W$. We now describe a way to turn Z into an F -vector space. Suppose we are given two functions $g, h \in Z$. We can define a new function $g + h$ mapping S into W by requiring that $g + h$ map each $x \in S$ to $g(x) + h(x)$ in W . In symbols, our definition says that $(g + h)(x) = g(x) + h(x)$ for all $x \in S$. Note that the plus symbol on the right side is the addition in the given vector space W , while the plus symbol on the left side is the new addition being defined on the set Z . It is routine to check that this new addition on Z (called *pointwise* addition of functions) satisfies the additive axioms for a vector space (see Table 1.2). In particular, the zero element of Z is the zero function that sends every $x \in S$ to 0_W , and the additive inverse of $f \in Z$ is the function that maps x to $-f(x)$ for all $x \in S$.

Now suppose we are given $c \in F$ and $g \in Z$. We define a new function $c \cdot g \in Z$ by

requiring that $(c \cdot g)(x) = c \cdot (g(x))$ for all $x \in S$. Note that the multiplication symbol on the right side is the given scalar multiplication of the F -vector space W , whereas the product symbol on the left side is the new scalar multiplication being defined on Z . One routinely verifies the remaining axioms for Z to be an F -vector space (see Table 1.4 and Exercise 1).

Here are some special cases of this construction. If $S = \{1, 2, \dots, n\}$ and W is the field F (which is an F -vector space), then Z is the vector space F^n of all n -tuples of elements of F (cf. Chapter 4). If $S = \{1, 2, \dots, m\} \times \{1, 2, \dots, n\}$ and $W = F$, then Z is the vector space of all $m \times n$ matrices with entries in F . Finally, suppose the set S is itself an F -vector space V . Then we conclude that the set of all functions from V to W is an F -vector space under pointwise operations on functions.

Fix F -vector spaces V and W . We define $\text{Hom}_F(V, W)$ to be the set of all F -linear transformations (vector space homomorphisms) from V to W . A function $f : V \rightarrow W$ belongs to $\text{Hom}_F(V, W)$ iff $f(x + y) = f(x) + f(y)$ and $f(cx) = cf(x)$ for all $x, y \in V$ and all $c \in F$. We claim that $\text{Hom}_F(V, W)$ is an F -vector space under the pointwise operations on functions defined above. It suffices to show that $\text{Hom}_F(V, W)$ is a subspace of the vector space of all functions from V to W . This follows since the zero function is F -linear, and sums and scalar multiples of F -linear functions are again F -linear (Exercise 2).

Next, we define the *dual space* of an F -vector space V to be the vector space $V^* = \text{Hom}_F(V, F)$ of all F -linear transformations from V to F . Here we are viewing the field F as a one-dimensional F -vector space where vector addition is addition in F , and scalar multiplication is the same as multiplication in the field F . Elements of V^* are often called *linear functionals*.

13.2 Dual Bases

Suppose V is a finite-dimensional F -vector space with ordered basis $B = (x_1, \dots, x_n)$. The basis B of V possesses the following *universal mapping property* (UMP). For any F -vector space W and any ordered list (y_1, \dots, y_n) of elements of W , there exists a unique F -linear map $T : V \rightarrow W$ such that $T(x_i) = y_i$ for all i with $1 \leq i \leq n$. To see this, note that every $v \in V$ has a unique expansion as a linear combination of elements of B , say

$$v = c_1x_1 + c_2x_2 + \cdots + c_nx_n \quad (c_i \in F).$$

If T is to be an F -linear map sending x_i to y_i , then T must send v to $c_1y_1 + \cdots + c_ny_n$. This observation proves the *uniqueness* of T . To show *existence* of T , define $T(\sum_{i=1}^n c_i x_i) = \sum_{i=1}^n c_i y_i$ for any $c_i \in F$, and check that T is indeed F -linear and sends each x_i to y_i (Exercise 5). Note that T is a well-defined function because each $v \in V$ can be written in the form $\sum_{i=1}^n c_i x_i$ in exactly one way.

The universal mapping property can be expressed informally as follows. To define a *linear map* from V to W , we are free to send each basis element x_i to any element of W . Having chosen the images of the basis elements, the entire map $T : V \rightarrow W$ is uniquely determined by linearity. Another way to say this is that each *function* $g : \{x_1, \dots, x_n\} \rightarrow W$ “extends by linearity” to a unique *linear map* $T_g : V \rightarrow W$.

As above, let $B = (x_1, \dots, x_n)$ be an ordered basis for the F -vector space V . For each fixed i between 1 and n , consider the function $g_i : \{x_1, \dots, x_n\} \rightarrow F$ such that $g_i(x_i) = 1_F$ and $g_i(x_j) = 0_F$ for all $j \neq i$. This function g_i extends by linearity to give a linear map $f_i : V \rightarrow F$ satisfying

$$f_i(c_1x_1 + \cdots + c_i x_i + \cdots + c_n x_n) = c_i \quad (c_1, \dots, c_n \in F).$$

Each f_i is an element of $V^* = \text{Hom}_F(V, F)$. So, starting with the ordered basis B of V , we have constructed an ordered list $B^d = (f_1, \dots, f_n)$ of elements of V^* . The UMP shows that B^d is the *unique* ordered list such that, for all i, j with $1 \leq i, j \leq n$, $f_i(x_j) = \delta_{ij}$, where δ_{ij} is 1 for $i = j$ and 0 otherwise. We claim that B^d is actually an ordered basis of the F -vector space V^* , which is called the *dual basis to B*. We must check that B^d is a linearly independent list spanning V^* . To check linear independence, assume

$$d_1 f_1 + d_2 f_2 + \cdots + d_n f_n = 0 \quad (d_i \in F), \quad (13.1)$$

where 0 denotes the zero function in V^* , and the operations on the left side are pointwise operations on functions. We need to prove that every $d_i = 0_F$. To do this, fix an index i , and evaluate both sides of (13.1) at the element x_i . We obtain

$$d_1 f_1(x_i) + \cdots + d_i f_i(x_i) + \cdots + d_n f_n(x_i) = 0(x_i).$$

By the definition of the f_j 's, the left side reduces to d_i , whereas the right side is zero. Thus $d_i = 0$ for each i , so B^d is a linearly independent list.

Next, let us see that B^d spans V^* . Given an arbitrary $h \in V^*$, we must express h as an F -linear combination of elements of B^d . We claim that, in the space V^* , we have

$$h = h(x_1)f_1 + h(x_2)f_2 + \cdots + h(x_n)f_n. \quad (13.2)$$

(Observe that each $h(x_i)$ is a scalar, since h maps V to F .) Note that both sides of this equation are linear maps from V to F . Now, two *linear* functions with domain V are equal iff they have the same value at each basis element x_i appearing in B . (This fact can be checked directly, but it also follows from the uniqueness assertion in the universal mapping property.) So, let us see whether the two sides of (13.2) agree at a fixed basis element x_i . The left side sends x_i to $h(x_i)$, and the right side sends x_i to

$$h(x_1)f_1(x_i) + \cdots + h(x_i)f_i(x_i) + \cdots + h(x_n)f_n(x_i) = h(x_i).$$

Thus (13.2) does hold, and B^d does span V .

Every vector space V has at least one basis, and all bases of V have the same cardinality, which is the dimension of V . Since the list B^d has the same length as the list B , the dual basis construction proves that *if V is a finite-dimensional vector space over F , then $\dim(V^*) = \dim(V)$* .

13.3 Zero Sets

Let V be an n -dimensional vector space over the field F . We are now ready to set up the correspondences between subsets of V^* and subsets of V that were mentioned in the introduction to this chapter. First we describe the map from subsets of V^* to subsets of V . Let S be any subset of V^* , so S is some collection of linear functionals mapping V into F . Define the *zero set* of S to be

$$\mathcal{Z}(S) = \{x \in V : f(x) = 0_F \text{ for all } f \in S\}.$$

When $V = F^n$, we can think of $\mathcal{Z}(S)$ as the solution set of the system of homogeneous linear equations $f(v) = 0$ as f ranges over S (cf. Exercise 15).

Here are some properties of the zero-set operator \mathcal{Z} :

- For all subsets S of V^* , $\mathcal{Z}(S)$ is a subspace of V . To prove this, first note that $0_V \in \mathcal{Z}(S)$ because all maps $f \in S$ (being linear) must send 0_V to 0_F . Second, assume $x, y \in \mathcal{Z}(S)$. For any $f \in S$, $f(x+y) = f(x) + f(y) = 0 + 0 = 0$, so $x+y \in \mathcal{Z}(S)$. Third, assume $x \in \mathcal{Z}(S)$ and $c \in F$. For any $f \in S$, $f(cx) = cf(x) = c0 = 0$, so $cx \in \mathcal{Z}(S)$.
 - \mathcal{Z} is inclusion-reversing, i.e., $S \subseteq T$ in V^* implies $\mathcal{Z}(T) \subseteq \mathcal{Z}(S)$ in V . To prove this, assume $S \subseteq T$ and $x \in \mathcal{Z}(T)$. For all $f \in S$, f is also in T , and so $f(x) = 0$. This proves that $x \in \mathcal{Z}(S)$.
 - $\mathcal{Z}(\emptyset) = V$. By definition, $\mathcal{Z}(\emptyset)$ is a subset of V . If it were a proper subset, there would exist $x \in V$ with $x \notin \mathcal{Z}(\emptyset)$. In turn, this would imply the existence of $f \in \emptyset$ with $f(x) \neq 0$. This is impossible, since \emptyset has no members.
 - $\mathcal{Z}(V^*) = \{0_V\}$. Note that $0_V \in \mathcal{Z}(V^*)$, since $\mathcal{Z}(V^*)$ is a subspace of V . Given $v \neq 0$ in V , let us show that $v \notin \mathcal{Z}(V^*)$. We must find an $f \in V^*$ with $f(v) \neq 0$. For this, extend the one-element list (v) to a basis $B = (v, v_2, \dots, v_n)$ of V . The first function in the dual basis B^d is an element of V^* sending v to $1_F \neq 0_F$.
 - If S is any subset of V^* and $W = \langle S \rangle$ is the subspace of V^* generated by S , then $\mathcal{Z}(S) = \mathcal{Z}(W)$. Since $S \subseteq W$, we already know $\mathcal{Z}(W) \subseteq \mathcal{Z}(S)$. Next, fix $x \in \mathcal{Z}(S)$. To prove $x \in \mathcal{Z}(W)$, fix $f \in W$. We can write $f = d_1g_1 + \dots + d_kg_k$ for some $d_i \in F$ and $g_i \in S$. Evaluating at x gives $f(x) = d_1g_1(x) + \dots + d_kg_k(x) = d_10 + \dots + d_k0 = 0$. So $x \in \mathcal{Z}(W)$.
-

13.4 Annihilators

Next we describe the map from subsets of V to subsets of V^* . Let T be any subset of V , so T is some collection of points in the vector space V . Define the *annihilator* of T to be

$$\mathcal{A}(T) = \{g \in V^* : g(x) = 0 \text{ for all } x \in T\}.$$

Here are some properties of the annihilator operator \mathcal{A} :

- For all subsets T of V , $\mathcal{A}(T)$ is a subspace of V^* . We check this as follows. First, 0_{V^*} sends every element of V to 0_F . In particular, 0_{V^*} sends every element of T to 0_F , so $0_{V^*} \in \mathcal{A}(T)$. Second, assume $g, h \in \mathcal{A}(T)$. For any $x \in T$, $(g+h)(x) = g(x) + h(x) = 0 + 0 = 0$, so $g+h \in \mathcal{A}(T)$. Third, assume $g \in \mathcal{A}(T)$ and $c \in F$. For any $x \in T$, $(cg)(x) = c(g(x)) = c0 = 0$, so $cg \in \mathcal{A}(T)$.
- \mathcal{A} is inclusion-reversing, i.e., $S \subseteq T$ in V implies $\mathcal{A}(T) \subseteq \mathcal{A}(S)$ in V^* . To prove this, assume $S \subseteq T$ and $g \in \mathcal{A}(T)$. For all $x \in S$, x is also in T , and so $g(x) = 0$. This proves that $g \in \mathcal{A}(S)$.
- $\mathcal{A}(\emptyset) = V^*$. The proof is analogous to the proof that $\mathcal{Z}(\emptyset) = V$ (Exercise 18).
- $\mathcal{A}(V) = \{0_{V^*}\}$. This holds since a linear functional $f : V \rightarrow F$ belongs to $\mathcal{A}(V)$ iff $f(x) = 0$ for all x in the domain V iff f is the zero function on V iff f is the zero vector in the space V^* .
- If T is any subset of V and $W = \langle T \rangle$ is the subspace of V generated by T , then $\mathcal{A}(T) = \mathcal{A}(W)$. Since $T \subseteq W$, we already know $\mathcal{A}(W) \subseteq \mathcal{A}(T)$. Next, fix $f \in \mathcal{A}(T)$. To prove $f \in \mathcal{A}(W)$, fix $x \in W$. We can write $x = c_1x_1 + \dots + c_kx_k$ for some $c_i \in F$ and $x_i \in T$. Now, linearity of f gives $f(x) = \sum_{i=1}^k c_i f(x_i) = \sum_{i=1}^k c_i 0 = 0$. So $f \in \mathcal{A}(W)$.

13.5 Double Dual V^{**}

The reader has probably noticed the striking similarity between the properties of \mathcal{Z} and the properties of \mathcal{A} . This similarity is not accidental. To explain it, we need the notion of the *double dual space* for V . The double dual of V is the dual space of the vector space V^* , namely

$$V^{**} = (V^*)^* = \text{Hom}_F(V^*, F) = \text{Hom}_F(\text{Hom}_F(V, F), F).$$

For finite-dimensional V , we know $\dim(V) = \dim(V^*) = \dim(V^{**})$.

Let us spell out the definition of V^{**} in more detail. An element of V^{**} is an F -linear function $E : V^* \rightarrow F$. The function E takes as input another F -linear function $g : V \rightarrow F$ and produces as output a scalar $E(g) \in F$. The F -linearity of E means that $E(g + h) = E(g) + E(h)$ and $E(dg) = dE(g)$ for all $g, h \in V^*$ and $d \in F$.

Elements of V^{**} may seem very abstract and difficult to visualize, since these elements are functions that operate on other functions. However, we can construct some concrete examples of elements of V^{**} as follows. Suppose x is a fixed vector in V . Define a map $E_x : V^* \rightarrow F$ by letting $E_x(g) = g(x)$ for all $g \in V^*$. In other words, E_x takes as input a linear functional $g : V \rightarrow F$ and returns as output the value of g at the fixed point $x \in V$, namely $g(x)$. For this reason, the function E_x is called *evaluation at x* . To check that E_x really is F -linear, use the definitions of the pointwise operations in V^* to compute

$$E_x(g + h) = (g + h)(x) = g(x) + h(x) = E_x(g) + E_x(h) \text{ for all } g, h \in V^*;$$

$$E_x(dg) = (dg)(x) = d(g(x)) = d(E_x(g)) \text{ for all } g \in V^*, d \in F.$$

We now know that every E_x really is an element of V^{**} .

Let us define a map $\text{ev} : V \rightarrow V^{**}$ by setting $\text{ev}(x) = E_x$ for all $x \in V$. In words, ev sends each vector x in V to the function “evaluation at x ,” which belongs to the double dual of V . The crucial fact about ev is that ev is a one-to-one F -linear map, which is bijective for finite-dimensional V .

We must check that $\text{ev} : V \rightarrow V^{**}$ is an F -linear map. (This is not the same as checking that $\text{ev}(x) = E_x$ is F -linear for each $x \in V$, which we have already done.) To check linearity of ev , first fix $x, y \in V$, and ask whether $\text{ev}(x + y) = \text{ev}(x) + \text{ev}(y)$. In other words, is $E_{x+y} = E_x + E_y$ in V^{**} ? These two functions from V^* to F are indeed equal, since for any $g \in V^*$,

$$E_{x+y}(g) = g(x + y) = g(x) + g(y) = E_x(g) + E_y(g) = (E_x + E_y)(g).$$

Next, fix $x \in V$ and $c \in F$, and ask whether $\text{ev}(cx) = c\text{ev}(x)$. In other words, is $E_{cx} = cE_x$ in V^{**} ? The answer is affirmative, since for any $g \in V^*$,

$$E_{cx}(g) = g(cx) = cg(x) = c(E_x(g)) = (cE_x)(g).$$

Now that we know ev is linear, we can show ev is one-to-one by proving that for all $x \in V$, $\text{ev}(x) = 0_{V^{**}}$ implies $x = 0_V$. We prove the contrapositive. Fix a nonzero $x \in V$. As shown in the proof of $\mathcal{Z}(V^*) = \{0\}$ (which readily extends to the case of infinite-dimensional spaces — cf. Exercise 10), there exists $f \in V^*$ with $f(x) \neq 0$. Then $E_x(f) \neq 0$, so $E_x = \text{ev}(x)$ is not the zero function from V^* to F . Hence, $\text{ev}(x) \neq 0$.

We now have an injective linear map $\text{ev} : V \rightarrow V^{**}$. If V is finite-dimensional, then $\dim(V^{**}) = \dim(V)$ as remarked earlier. So the injective linear map ev must also be surjective since the domain and codomain have the same finite dimension.

Thus, for each finite-dimensional F -vector space V , we have a natural vector space

isomorphism $\text{ev} : V \rightarrow V^{**}$. The word *natural* is a technical mathematical adjective that will be explained later. The fact that ev is a bijection means that for every $E \in V^{**}$, there exists a unique $x \in V$ such that $E = E_x$ (evaluation at x). Linearity of ev means $E_{x+y} = E_x + E_y$ and $E_{cx} = cE_x$ for all $x, y \in V$ and $c \in F$, as we confirmed earlier.

For the rest of this section, assume V is finite-dimensional. We often use the isomorphism ev to identify V^{**} with V , thereby blurring the distinction between an element $x \in V$ and the associated evaluation map $E_x \in V^{**}$. This identification may seem confusing at first, but it actually clarifies certain issues and makes certain proofs less redundant. For instance, consider the relation between \mathcal{Z} and \mathcal{A} mentioned at the beginning of this subsection. Let \mathcal{A}^* be the annihilator operator for the vector space V^* , which sends subsets of V^* to subsets of V^{**} . Given a subset S of V^* , the definition of \mathcal{A}^* gives

$$\mathcal{A}^*(S) = \{E \in V^{**} : E(h) = 0 \text{ for all } h \in S\}.$$

But every $E \in V^{**}$ has the form E_x for a unique $x \in V$. Then $\text{ev}(x) = E_x \in \mathcal{A}^*(S)$ iff $E_x(h) = 0$ for all $h \in S$ iff $h(x) = 0$ for all $h \in S$ iff $x \in \mathcal{Z}(S)$. This proves that $\mathcal{A}^*(S)$ is the image of $\mathcal{Z}(S)$ under the isomorphism ev . Now, if we use ev to identify V with V^{**} , the previous statement says that $\mathcal{A}^*(S) = \mathcal{Z}(S)$.

Similarly, letting \mathcal{Z}^* be the zero-set operator for V^* (which maps subsets of V^{**} to subsets of V^*), we have $\mathcal{Z}^*(T) = \mathcal{A}(T)$ for all subsets T of V , under the identification of V^{**} with V . More precisely, we are asserting that $\mathcal{Z}^*(\text{ev}[T]) = \mathcal{A}(T)$ for all $T \subseteq V$. The reason is that $g \in \mathcal{A}(T)$ iff $g(x) = 0$ for all $x \in T$ iff $E_x(g) = 0$ for all $x \in T$ iff g is sent to zero by all functions in $\text{ev}[T]$ iff $g \in \mathcal{Z}^*(\text{ev}[T])$.

With these facts in hand, we can see why the properties of the operators \mathcal{Z} and \mathcal{A} exactly mirror each other. For example, given that the annihilator operator (on *any* vector space) has already been shown to be inclusion-reversing, it follows that the zero-set operator is also inclusion-reversing, because

$$S \subseteq T \subseteq V^* \Rightarrow \mathcal{Z}(T) = \mathcal{A}^*(T) \subseteq \mathcal{A}^*(S) = \mathcal{Z}(S).$$

As another example of this process, let us prove the two results $\mathcal{Z}(\mathcal{A}(T)) \supseteq T$ for all $T \subseteq V$ and $\mathcal{A}(\mathcal{Z}(S)) \supseteq S$ for all $S \subseteq V^*$. For the first result, assume T is a subset of V , fix $x \in T$, and prove $x \in \mathcal{Z}(\mathcal{A}(T))$. We must show $g(x) = 0$ for all $g \in \mathcal{A}(T)$. But this follows from the very definition of $\mathcal{A}(T)$, since $x \in T$. The first result therefore holds for *all* vector spaces V . In particular, applying this result to the vector space V^* instead of V , we see that $\mathcal{Z}^*(\mathcal{A}^*(S)) \supseteq S$ for all $S \subseteq V^*$. As seen above, this means $\mathcal{A}(\mathcal{Z}(S)) \supseteq S$ for all $S \subseteq V^*$.

As a final example of working with the double dual, we show that the dual basis construction works “in reverse.” Suppose $B^* = (f_1, \dots, f_n)$ is any ordered basis of V^* . We show there exists a unique dual basis $(B^*)^D = (x_1, \dots, x_n)$ of V such that $f_i(x_j) = \delta_{ij}$ for all $i, j \in [n]$. To get this basis, first form the dual basis (in the old sense) $(B^*)^d = (E_1, \dots, E_n)$, which is a basis of V^{**} . Then write $E_j = E_{x_j} = \text{ev}(x_j)$ for a unique $x_j \in V$. Since ev is an isomorphism, (x_1, \dots, x_n) is an ordered basis for V , and moreover $f_i(x_j) = E_{x_j}(f_i) = E_j(f_i) = \delta_{ij}$. The uniqueness of the basis (x_1, \dots, x_n) follows from the uniqueness of (E_1, \dots, E_n) and the bijectivity of ev . One can also show that, for any ordered basis B of V , $B^{dD} = B$. To see why, write $B = (x_1, \dots, x_n)$, $B^d = (f_1, \dots, f_n)$, and $B^{dD} = (z_1, \dots, z_n)$. We must have $x_i = z_i$ for each i , because $f_i(z_j) = \delta_{i,j}$, and we know B is the *unique* ordered basis of V satisfying this condition.

13.6 Correspondence between Subspaces of V and V^*

Let V be a finite-dimensional vector space over a field F . Let $\mathcal{P}(V)$ denote the set of all subsets of V , and let $\mathcal{S}(V)$ denote the set of all subspaces of V . Define $\mathcal{P}(V^*)$ and $\mathcal{S}(V^*)$ similarly. In the previous subsections, we introduced two inclusion-reversing maps

$$\mathcal{Z} : \mathcal{P}(V^*) \rightarrow \mathcal{S}(V), \quad \mathcal{A} : \mathcal{P}(V) \rightarrow \mathcal{S}(V^*).$$

These maps are *not* one-to-one. For instance, if W is a subspace of V that is generated by a proper subset S , then $S \neq W$ but $\mathcal{A}(S) = \mathcal{A}(W)$.

To rectify this problem, we restrict the domains of \mathcal{Z} and \mathcal{A} to the set of subspaces of V^* and V , respectively, giving two inclusion-reversing maps

$$\mathcal{Z} : \mathcal{S}(V^*) \rightarrow \mathcal{S}(V), \quad \mathcal{A} : \mathcal{S}(V) \rightarrow \mathcal{S}(V^*).$$

The key result is: *for all finite-dimensional F -vector spaces V , the maps \mathcal{Z} and \mathcal{A} are mutually inverse inclusion-reversing bijections between $\mathcal{S}(V^*)$ and $\mathcal{S}(V)$ satisfying*

$$\dim(W) + \dim(\mathcal{Z}(W)) = \dim(V) = \dim(Y) + \dim(\mathcal{A}(Y)) \quad \text{for all } W \in \mathcal{S}(V^*), Y \in \mathcal{S}(V). \quad (13.3)$$

Furthermore, the original (unrestricted) \mathcal{Z} and \mathcal{A} operators satisfy

$$\mathcal{A}(\mathcal{Z}(S)) = \langle S \rangle, \quad \mathcal{Z}(\mathcal{A}(T)) = \langle T \rangle \quad \text{for all } S \in \mathcal{P}(V^*), T \in \mathcal{P}(V).$$

Let us prove the dimension formula $\dim(V) = \dim(Y) + \dim(\mathcal{A}(Y))$, where Y is a fixed subspace of V . Let (y_1, \dots, y_k) be an ordered basis for Y . We can extend this to an ordered basis $B = (y_1, \dots, y_k, y_{k+1}, \dots, y_n)$ for V , where $n = \dim(V)$. Let $B^d = (f_1, \dots, f_k, f_{k+1}, \dots, f_n)$ be the dual basis for V^* . We will prove that (f_{k+1}, \dots, f_n) is a basis for $\mathcal{A}(Y)$, which will imply the needed relation $\dim(\mathcal{A}(Y)) = n - k = \dim(V) - \dim(Y)$. Now (f_{k+1}, \dots, f_n) is a linearly independent list in V^* , being a sublist of the ordered basis B^d . A generic element $f \in V^*$ can be written uniquely as

$$f = d_1 f_1 + \cdots + d_k f_k + d_{k+1} f_{k+1} + \cdots + d_n f_n \quad (d_j \in F).$$

We claim $f \in \mathcal{A}(Y)$ iff $d_1 = \cdots = d_k = 0$. On one hand, $f \in \mathcal{A}(Y)$ iff $f \in \mathcal{A}(\{y_1, \dots, y_k\})$, since these vectors generate the subspace Y . On the other hand, $f(y_i) = \sum_{j=1}^n d_j f_j(y_i) = d_i$ for all i . Therefore, $f \in \mathcal{A}(Y)$ iff $f(y_1) = \cdots = f(y_k) = 0$ iff $d_1 = \cdots = d_k = 0$. It is now evident that $\mathcal{A}(Y)$ consists precisely of all F -linear combinations of the linearly independent vectors f_{k+1}, \dots, f_n . Thus these vectors are indeed a basis for $\mathcal{A}(Y)$.

The proof of (13.3) can now be completed quickly. Applying the result of the last paragraph to V^* instead of V , and using the identification of V^{**} with V , we get

$$\dim(V) = \dim(V^*) = \dim(W) + \dim(\mathcal{A}^*(W)) = \dim(W) + \dim(\mathcal{Z}(W)) \quad \text{for all } W \in \mathcal{S}(V^*).$$

Next, we will show that the restricted \mathcal{Z} and \mathcal{A} operators are two-sided inverses of each other, hence both are bijections. Suppose Y is any subspace of V . We have seen that $\mathcal{Z}(\mathcal{A}(Y))$ is a subspace of V containing Y . But on the other hand, the dimension formulas in (13.3) show that

$$\dim(\mathcal{Z}(\mathcal{A}(Y))) = \dim(V) - \dim(\mathcal{A}(Y)) = \dim(V) - (\dim(V) - \dim(Y)) = \dim(Y).$$

Since the dimensions are equal and finite, $\mathcal{Z}(\mathcal{A}(Y))$ cannot properly contain Y . So

$\mathcal{Z}(\mathcal{A}(Y)) = Y$ holds. The same argument establishes that $\mathcal{A}(\mathcal{Z}(W)) = W$ for all $W \in \mathcal{S}(V^*)$. Finally, if S is any subset of V^* , we compute

$$\mathcal{A}(\mathcal{Z}(S)) = \mathcal{A}(\mathcal{Z}(\langle S \rangle)) = \langle S \rangle,$$

and similarly $\mathcal{Z}(\mathcal{A}(T)) = \langle T \rangle$ for all $T \subseteq V$.

We pause to recall some definitions from the theory of posets (see the Appendix for more details). A *poset* (partially ordered set) consists of a set Z and a relation \leq on Z that is *reflexive* (for all $x \in Z$, $x \leq x$); *antisymmetric* (for all $x, y \in Z$, if $x \leq y$ and $y \leq x$ then $x = y$); and *transitive* (for all $x, y, z \in Z$, if $x \leq y$ and $y \leq z$ then $x \leq z$). Given $S \subseteq Z$, an *upper bound* for S is an element $z \in Z$ such that for all $x \in S$, $x \leq z$. A *least upper bound* for S , denoted $\sup S$, is an upper bound w for S such that $w \leq z$ for all upper bounds z for S . *Lower bounds* for S and the *greatest lower bound* for S (denoted $\inf S$) are defined similarly. Subsets of posets do not always have least upper bounds or greatest lower bounds, but $\sup S$ and $\inf S$ are unique when they exist. The poset Z is called a *lattice* iff every two-element subset of Z has a least upper bound and a greatest lower bound. Z is a *complete lattice* iff every nonempty subset of Z has a least upper bound and a greatest lower bound.

Observe that $\mathcal{S}(V)$ and $\mathcal{S}(V^*)$ are both posets ordered by set inclusion: $X \leq Y$ means $X \subseteq Y$. One may routinely check that if X and Y are subspaces of V , then in the poset $\mathcal{S}(V)$, $\inf\{X, Y\} = X \cap Y$ and $\sup\{X, Y\} = X + Y = \{x + y : x \in X, y \in Y\}$. So $\mathcal{S}(V)$ is a lattice. More generally, if $\{X_i : i \in I\}$ is any indexed collection of subspaces of V , then $\inf_{i \in I} X_i = \bigcap_{i \in I} X_i$ and $\sup_{i \in I} X_i = \sum_{i \in I} X_i$, where $\sum_{i \in I} X_i$ is the set of all finite sums $x_{i_1} + \cdots + x_{i_k}$ with $x_{i_j} \in X_{i_j}$. So $\mathcal{S}(V)$ is a complete lattice; similarly, $\mathcal{S}(V^*)$ is a complete lattice. We have shown that the maps $\mathcal{Z} : \mathcal{S}(V^*) \rightarrow \mathcal{S}(V)$ and $\mathcal{A} : \mathcal{S}(V) \rightarrow \mathcal{S}(V^*)$ are bijections that reverse inclusions; in the theory of posets, we would therefore say that these maps are *poset anti-isomorphisms*. It follows from this fact and the definitions of sup and inf that \mathcal{Z} and \mathcal{A} interchange least upper bounds and greatest lower bounds (Exercise 37). In other words,

$$\mathcal{Z}(X + Y) = \mathcal{Z}(X) \cap \mathcal{Z}(Y), \quad \mathcal{Z}(X \cap Y) = \mathcal{Z}(X) + \mathcal{Z}(Y) \quad \text{for } X, Y \in \mathcal{S}(V^*); \quad (13.4)$$

$$\mathcal{A}(X + Y) = \mathcal{A}(X) \cap \mathcal{A}(Y), \quad \mathcal{A}(X \cap Y) = \mathcal{A}(X) + \mathcal{A}(Y) \quad \text{for } X, Y \in \mathcal{S}(V); \quad (13.5)$$

and similar formulas hold for the sup and inf of indexed collections of subspaces.

13.7 Dual Maps

Suppose V and W are vector spaces over a field F , and $T : V \rightarrow W$ is a linear transformation. If $f : W \rightarrow F$ is any element of W^* , observe that the composite map $f \circ T$ is a linear map from V to F , i.e., an element of V^* . Thus, we can define a function $T^* : W^* \rightarrow V^*$ by setting $T^*(f) = f \circ T$ for all $f \in W^*$. The function T^* is itself a linear map, since

$$T^*(f + g) = (f + g) \circ T = (f \circ T) + (g \circ T) = T^*(f) + T^*(g) \quad \text{for all } f, g \in W^*;$$

$$T^*(cf) = (cf) \circ T = c(f \circ T) = cT^*(f) \quad \text{for all } c \in F, f \in W^*.$$

(To verify the equality $(f + g) \circ T = (f \circ T) + (g \circ T)$, note that both sides are functions from V to F sending each $x \in V$ to $f(T(x)) + g(T(x))$. A similar calculation justifies the equality $(cf) \circ T = c(f \circ T)$.) We call the linear map $T^* : W^* \rightarrow V^*$ the *dual map* to $T : V \rightarrow W$.

Let us note some algebraic properties of this construction. First, the passage from T to T^* is *linear*. This means that if $T, U : V \rightarrow W$ are linear maps and $c \in F$ is a scalar, then $(T+U)^* = T^* + U^*$ and $(cT)^* = c(T^*)$. For example, the first equality holds because both sides are functions from W^* to V^* sending $g \in W^*$ to $g \circ (T+U) = (g \circ T) + (g \circ U)$. Similarly, both $(cT)^*$ and $c(T^*)$ send g to $c(g \circ T)$. Also, the passage from T to T^* is *functorial*. This means that if $T : V \rightarrow W$ and $S : W \rightarrow Z$ are linear maps, then $(S \circ T)^* = T^* \circ S^*$; and if $\text{id} : V \rightarrow V$ is the identity map, then id^* is the identity map on V^* . To check the first assertion, note

$$(S \circ T)^*(h) = h \circ (S \circ T) = (h \circ S) \circ T = (S^*(h)) \circ T = T^*(S^*(h)) = (T^* \circ S^*)(h) \quad \text{for } h \in Z^*.$$

The second assertion follows since $\text{id}^*(f) = f \circ \text{id} = f$ for $f \in V^*$.

Suppose $B = (v_1, \dots, v_n)$ is an ordered basis for V , and $C = (w_1, \dots, w_m)$ is an ordered basis for W . The matrix of T relative to the bases B and C is the unique $m \times n$ matrix of scalars $A = (a_{ij})$ such that

$$T(v_j) = \sum_{i=1}^m a_{ij} w_i \quad (a_{ij} \in F, 1 \leq j \leq n).$$

We have seen that the dual basis $B^d = (f_1, \dots, f_n)$ is an ordered basis for V^* , and $C^d = (g_1, \dots, g_m)$ is an ordered basis for W^* . What is the matrix of the linear map T^* relative to the bases C^d and B^d ? By definition, this is the unique $n \times m$ matrix of scalars $A' = (a'_{rs})$ such that

$$T^*(g_s) = \sum_{r=1}^n a'_{rs} f_r \quad (a'_{rs} \in F, 1 \leq s \leq m).$$

We claim that these equations hold if we take $a'_{rs} = a_{sr}$, so that A' is the *transpose* of the matrix A . To see this, fix s with $1 \leq s \leq m$, and let us check whether the two functions $T^*(g_s)$ and $\sum_{r=1}^n a_{sr} f_r$ are equal. Both functions are linear maps from V to F , so it suffices to see that they take the same value at each basis element v_k , where $1 \leq k \leq n$. On one hand,

$$[T^*(g_s)](v_k) = (g_s \circ T)(v_k) = g_s(T(v_k)) = g_s \left(\sum_{i=1}^m a_{ik} w_i \right) = a_{sk}.$$

The last equality uses the fact that C^d is the dual basis for C . On the other hand,

$$\left[\sum_{r=1}^n a_{sr} f_r \right] (v_k) = \sum_{r=1}^n a_{sr} f_r(v_k) = a_{sk}$$

since B^d is the dual basis for B . To summarize, we have shown that *if the matrix A represents a linear map T relative to certain bases, then the transpose of A represents the dual map T^* relative to the dual bases*. For this reason, the map T^* is sometimes called the *transpose* of the map T .

Given a linear map $T : V \rightarrow W$, we have the dual map $T^* : W^* \rightarrow V^*$. Iterating the construction produces the *double dual map* $T^{**} : V^{**} \rightarrow W^{**}$, which is defined by $T^{**}(E) = E \circ T^*$ for $E \in V^{**}$. Expanding this definition further, we see that $T^{**}(E)$ maps $g \in W^*$ to $E(T^*(g)) = E(g \circ T)$.

Now we can explain precisely what it means to say that the map $\text{ev} = \text{ev}_V : V \rightarrow V^{**}$ is *natural*. Naturality means that for any vector spaces V and W and any linear map

$T : V \rightarrow W$, the following diagram commutes:

$$\begin{array}{ccc} V & \xrightarrow{T} & W \\ \text{ev}_V \downarrow & & \downarrow \text{ev}_W \\ V^{**} & \xrightarrow{T^{**}} & W^{**} \end{array}$$

(i.e., $\text{ev}_W \circ T = T^{**} \circ \text{ev}_V$). To check this, fix $x \in V$, and ask whether the two functions $\text{ev}_W(T(x)) = E_{T(x)}$ and $T^{**}(\text{ev}_V(x)) = T^{**}(E_x)$ in W^{**} are equal. For this, fix $g \in W^*$, and apply each function to g . The first function produces $E_{T(x)}(g) = g(T(x))$. The second function gives

$$[T^{**}(E_x)](g) = (E_x \circ T^*)(g) = E_x(T^*(g)) = E_x(g \circ T) = (g \circ T)(x) = g(T(x)).$$

Thus, the functions are equal, and the diagram does commute.

In the finite-dimensional case, we say that V^{**} is *naturally* isomorphic to V under the map ev_V , since we can “lift” linear maps $T : V \rightarrow W$ to linear maps $T^{**} : V^{**} \rightarrow W^{**}$. On the other hand, even though V and V^* are isomorphic, it can be shown that there is no *natural* isomorphism between V and V^* when $\dim(V) > 1$ (Exercise 40). As we will soon see, this situation can be partially remedied if V has appropriate additional structure (e.g., if V is an inner product space).

13.8 Nondegenerate Bilinear Forms

Let V be a vector space over the field F . A function $B : V \times V \rightarrow F$ is called a *bilinear form on V* iff

$$B(x_1 + x_2, y) = B(x_1, y) + B(x_2, y), \quad B(cx, y) = cB(x, y) \quad \text{for all } c \in F, x, y, x_1, x_2 \in V;$$

$$B(x, y_1 + y_2) = B(x, y_1) + B(x, y_2), \quad B(x, cy) = cB(x, y) \quad \text{for all } c \in F, x, y, y_1, y_2 \in V.$$

The first two conditions say that, for each fixed $y \in V$, the function sending each $x \in V$ to $B(x, y) \in F$ is a *linear* map from V to F , i.e., an element of V^* . We sometimes denote this element of V^* by the notation $B(\cdot, y)$. Similarly, the second two conditions state that, for each fixed $x \in V$, the function $B(x, \cdot)$ mapping each $y \in V$ to $B(x, y) \in F$ is an element of V^* .

We can use these observations to define two maps $L_B, R_B : V \rightarrow V^*$. Specifically, for all $x, y \in V$, set

$$L_B(x) = B(x, \cdot) \text{ and } R_B(y) = B(\cdot, y).$$

The comments above show that L_B and R_B do map into the stated codomain V^* . Furthermore, each of the maps L_B and R_B are *linear*. For instance, $L_B(x+z) = L_B(x) + L_B(z)$ for all $x, z \in V$ because, for each $y \in V$,

$$L_B(x+z)(y) = B(x+z, y) = B(x, y) + B(z, y) = L_B(x)(y) + L_B(z)(y) = (L_B(x) + L_B(z))(y).$$

We similarly verify that $L_B(cx) = cL_B(x)$ (where $x \in V$ and $c \in F$) and that R_B is linear.

Assume V is finite-dimensional. One can now ask whether the linear maps L_B and R_B are *isomorphisms* from V to V^* . Since $\dim(V) = \dim(V^*)$, a linear map from V to V^* is an isomorphism iff the map is one-to-one, which holds iff the kernel of the map consists

of zero alone. We say that the bilinear form B is *nondegenerate on the left* iff L_B is an isomorphism, and B is called *nondegenerate on the right* iff R_B is an isomorphism. By the previous remark, B is nondegenerate on the left iff for all $x \in V$, $(B(x, y) = 0 \text{ for all } y \in V)$ implies $x = 0$; and B is nondegenerate on the right iff for all $y \in V$, $(B(x, y) = 0 \text{ for all } x \in V)$ implies $y = 0$. One can show that B is nondegenerate on the left iff B is nondegenerate on the right iff a certain matrix representing B is invertible (Exercise 51). So we can speak of *nondegenerate* bilinear forms.

Let us consider the particular case where B is a *symmetric* bilinear form, which means $B(x, y) = B(y, x)$ for all $x, y \in V$. If B is symmetric and nondegenerate, then $L_B : V \rightarrow V^*$ is a vector space isomorphism. This means that *for every linear map $f : V \rightarrow F$, there exists a unique vector $x \in V$ such that $f(y) = B(x, y)$ for all $y \in V$* . Existence of x is surjectivity of L_B ; uniqueness of x is injectivity of L_B ; and we have $f = L_B(x)$. In this situation, we have a convenient concrete description of the elements of V^* : each element of V^* is “multiplication by x on the left” (relative to B) for some uniquely determined $x \in V$. Compare this to the situation for V^{**} , where every element of V^{**} was given by *evaluation at* some uniquely determined $x \in V$. If B is fixed, we can use the isomorphism $L_B = R_B$ to identify V with V^* , just as we used ev earlier to identify V with V^{**} .

As an application of these ideas, let B be a symmetric, nondegenerate bilinear form on a finite-dimensional vector space V . We will show that *for each linear operator $T : V \rightarrow V$, there exists a unique linear operator $T' : V \rightarrow V$ satisfying*

$$B(T(x), y) = B(x, T'(y)) \text{ for all } x, y \in V. \quad (13.6)$$

T' is called the *adjoint* of T relative to B .

To show this, fix $y \in V$. First note that the map sending each $x \in V$ to $B(T(x), y) \in F$ is none other than the map $R_B(y) \circ T = T^*(R_B(y)) \in V^*$. On the other hand, the map sending each $x \in V$ to $B(x, T'(y)) \in F$ is $R_B(T'(y))$. So the required condition will hold for all x and for this fixed y iff

$$(T^* \circ R_B)(y) = (R_B \circ T')(y).$$

Therefore, the required condition will hold for all x and all y in V iff $T^* \circ R_B = R_B \circ T'$ iff this diagram commutes:

$$\begin{array}{ccc} V & \xrightarrow{R_B} & V^* \\ T' \downarrow & & \downarrow T^* \\ V & \xrightarrow{R_B} & V^* \end{array}$$

Since R_B is an isomorphism, there is a unique linear map $T' : V \rightarrow V$ with the required properties, namely $T' = R_B^{-1} \circ T^* \circ R_B$. If we use the isomorphism R_B to identify V with V^* , we see that the adjoint map $T' : V \rightarrow V$ is essentially the same as the dual map $T^* : V^* \rightarrow V^*$.

13.9 Real Inner Product Spaces

We can use the ideas in the previous section to derive some fundamental properties of inner product spaces. A *real inner product space* is a vector space V over \mathbb{R} together with a symmetric bilinear form B on V (typically denoted $B(x, y) = \langle x, y \rangle$ and called the *scalar product* on V) such that $B(x, x) > 0$ for all nonzero $x \in V$. Note that $B(0, 0) = 0$ by

bilinearity. Such a scalar product is automatically nondegenerate, because $B(x, y) = 0$ for all $y \in V$ implies $B(x, x) = 0$, which implies $x = 0$. So we can apply all our previous results on nondegenerate symmetric bilinear forms.

In particular, for finite-dimensional V , every linear functional $f : V \rightarrow \mathbb{R}$ has the form $f(y) = \langle y, x \rangle$ for a uniquely determined $x \in V$. Similarly, for any linear operator T on V , there is a unique *adjoint operator* T' on V characterized by

$$\langle T(x), y \rangle = \langle x, T'(y) \rangle \quad \text{for all } x, y \in V.$$

Suppose $C = (x_1, \dots, x_n)$ is an ordered basis for V . We claim there exists a unique *dual basis* $C' = (y_1, \dots, y_n)$ for V such that $\langle x_i, y_j \rangle = \delta_{ij}$ for all $i, j \in [n]$. To see this, let $C^d = (f_1, \dots, f_n)$ be the unique dual basis of C in V^* , and then let y_1, \dots, y_n be the unique elements of V such that f_j is right-multiplication by y_j . Note that (y_1, \dots, y_n) is an ordered basis of V , since it is the image of an ordered basis under the isomorphism R_B^{-1} .

Finally, let us discuss orthogonal complements. Let W be any subspace of V . Say $\dim(W) = k$ and $\dim(V) = n$. We know $\mathcal{A}(W) \subseteq V^*$ is an $(n - k)$ -dimensional subspace consisting of linear functionals that annihilate W . Taking the image of $\mathcal{A}(W)$ under the isomorphism $R_B^{-1} : V^* \rightarrow V$, we obtain the set

$$W^\perp = \{y \in V : \langle x, y \rangle = 0 \text{ for all } x \in W\},$$

which is called the *orthogonal complement of W in V* . Since this set is the image of $\mathcal{A}(W)$ under an isomorphism, W^\perp is in fact a subspace such that

$$\dim(W) + \dim(W^\perp) = n.$$

Similarly, if Y is a subspace of V^* and we identify V^* with V using R_B , we see that

$$\mathcal{Z}(Y) = \{x \in V : \langle x, y \rangle = 0 \text{ for all } y \in Y\} = Y^\perp.$$

Thus, the \mathcal{A} and \mathcal{Z} operators from our discussion of dual spaces collapse into a single “orthogonality operator” on V . This operator, which maps each subspace W of V to W^\perp , is a poset anti-isomorphism of $\mathcal{S}(V)$ such that $\dim(W) + \dim(W^\perp) = \dim(V)$. Moreover, since $\mathcal{Z}(\mathcal{A}(W)) = W$ for any subspace W , we have $W^{\perp\perp} = W$. More generally, for any subset S of V , $S^{\perp\perp} = \langle S \rangle$, the subspace of V generated by S . If W is any subspace of V , then $W \cap W^\perp = \{0\}$ since $z \in W \cap W^\perp$ implies $\langle z, z \rangle = 0$, forcing $z = 0$. Combining this with the dimension formula $\dim(W) + \dim(W^\perp) = \dim(V)$, we see that $V = W \oplus W^\perp$ is the *direct sum* of any subspace W and its orthogonal complement. We warn the reader that the last two results ($W \cap W^\perp = \{0\}$ and $V = W \oplus W^\perp$) do *not* always hold for an arbitrary symmetric nondegenerate bilinear form, although the other results mentioned (in particular, $\dim(W) + \dim(W^\perp) = \dim(V)$) are still true.

13.10 Complex Inner Product Spaces

Let V be a finite-dimensional vector space over the field \mathbb{C} of complex numbers. A function $B : V \times V \rightarrow \mathbb{C}$ (typically denoted $B(x, y) = \langle x, y \rangle$ for $x, y \in V$) is called a *complex inner product* for V iff these four conditions hold:

$$\langle x_1 + x_2, y \rangle = \langle x_1, y \rangle + \langle x_2, y \rangle, \quad \langle cx, y \rangle = c\langle x, y \rangle,$$

$$\langle y, x \rangle = \overline{\langle x, y \rangle}, \quad \text{and} \quad (x \neq 0 \Rightarrow \langle x, x \rangle \in \mathbb{R}^+) \quad \text{for all } c \in \mathbb{C}, x, y, x_1, x_2 \in V.$$

The bar in the third equation denotes complex conjugation ($\overline{a+ib} = a - ib$ for $a, b \in \mathbb{R}$). It follows that $\langle x, x \rangle$ must be real for all $x \in V$ (since $\langle x, x \rangle = \overline{\langle x, x \rangle}$), and the fourth condition requires that this quantity be a strictly positive real number for $x \neq 0$. The first two formulas state that $R_B(y) = \langle \cdot, y \rangle$ is a \mathbb{C} -linear map for each fixed $y \in V$. Using the third formula, we find on the other hand that, for all $x, y, y_1, y_2 \in V$ and $c \in \mathbb{C}$,

$$\begin{aligned}\langle x, y_1 + y_2 \rangle &= \overline{\langle y_1 + y_2, x \rangle} = \overline{\langle y_1, x \rangle + \langle y_2, x \rangle} = \overline{\langle y_1, x \rangle} + \overline{\langle y_2, x \rangle} = \langle x, y_1 \rangle + \langle x, y_2 \rangle; \\ \langle x, cy \rangle &= \overline{\langle cy, x \rangle} = \overline{c\langle y, x \rangle} = \bar{c} \cdot \overline{\langle y, x \rangle} = \bar{c}\langle x, y \rangle.\end{aligned}$$

So, for each fixed $x \in V$, $L_B(x) = \langle x, \cdot \rangle$ is almost a linear map, except for the extra conjugation of complex scalars. Let us call a map $T : V \rightarrow W$ between two complex vector spaces *semi-linear* iff $T(x+y) = T(x) + T(y)$ and $T(cx) = \bar{c}T(x)$ for all $x, y \in V$ and all $c \in \mathbb{C}$. A bijective semi-linear map will be called a *semi-isomorphism*.

We list some facts about semi-linear maps whose routine verifications are left as exercises. A semi-linear map $T : V \rightarrow W$ is injective iff $\ker(T) = \{x \in V : T(x) = 0_W\}$ is $\{0_V\}$. The composition of two semi-linear maps is a linear map. The composition of a linear map and a semi-linear map is a semi-linear map. Semi-isomorphisms map linearly independent lists to linearly independent lists and bases to bases. The inverse of a semi-isomorphism is a semi-isomorphism. The image or inverse image of a subspace under a semi-linear map is again a subspace, which has the same dimension as the first subspace in the case of a semi-isomorphism. The rank-nullity formula, $\dim(V) = \dim(\ker(T)) + \dim(\text{img}(T))$, is valid for a semi-linear map $T : V \rightarrow W$. Let V^{*s} be the set of semi-linear maps from V to \mathbb{C} ; this is a complex vector space (under pointwise operations) with $\dim(V^{*s}) = \dim(V) < \infty$.

With this terminology in hand, we can now say that $L_B(x) = \langle x, \cdot \rangle$ is a semi-linear map for each fixed $x \in V$. Furthermore, one can check that the map L_B itself (which maps V into V^{*s}) is *linear*, whereas the map $R_B : V \rightarrow V^*$ is *semi-linear*. For example, given a fixed $x \in V$, semi-linearity of $L_B(x)$ says that $L_B(x)(cy) = \langle x, cy \rangle = \bar{c}\langle x, y \rangle = \bar{c}L_B(x)(y)$ for all $y \in V$ and all $c \in \mathbb{C}$; whereas linearity of L_B says that $L_B(cx) = \langle cx, \cdot \rangle = c\langle x, \cdot \rangle = cL_B(x)$ for all $x \in V$ and all $c \in \mathbb{C}$.

What are the kernels of L_B and R_B in this setup? Given $x \in V$, if $L_B(x) = 0$, then in particular $\langle x, x \rangle = 0$, forcing $x = 0$. Similarly, for all $x \in V$, $R_B(x) = 0$ implies $x = 0$. So L_B and R_B are both injective, hence bijective since $\dim(V^{*s}) = \dim(V) = \dim(V^*)$. Thus, L_B gives an isomorphism of V and V^{*s} , whereas R_B gives a semi-isomorphism of V and V^* . In other words, every (linear) $f \in V^*$ has the form $f = R_B(y) = \langle \cdot, y \rangle$ for a unique $y \in V$, whereas every (semi-linear) $g \in V^{*s}$ has the form $g = L_B(x) = \langle x, \cdot \rangle$ for a unique $x \in V$. This allows us to identify the vector spaces V and V^{*s} . We can almost identify the vector spaces V and V^* as well, except for the extra conjugations of scalars caused by the semi-isomorphism.

Now let W be any subspace of the complex inner product space V . We have the orthogonal complement

$$W^\perp = \{y \in V : \langle x, y \rangle = 0 \text{ for all } x \in W\}.$$

Note that, for $y \in V$, $y \in W^\perp$ iff $R_B(y) \in \mathcal{A}(W) \subseteq V^*$. This says that $\mathcal{A}(W)$ is the image of W^\perp under the semi-isomorphism R_B . Therefore, as before, W^\perp is a subspace of V with $\dim(W^\perp) = \dim(\mathcal{A}(W)) = \dim(V) - \dim(W)$. This dimension formula, together with $W \subseteq W^{\perp\perp}$, proves that $W = W^{\perp\perp}$. The map $W \mapsto W^\perp$ is an inclusion-reversing bijection on the lattice of subspaces of V . Finally, the condition $\langle x, x \rangle > 0$ for $x \neq 0$ forces $W \cap W^\perp = \{0\}$. So, as in the case of real inner product spaces, we have $V = W \oplus W^\perp$ for any subspace W .

The concept of dual bases carries over to the present setup as follows. Suppose

(x_1, \dots, x_n) is an ordered basis of the complex inner product space V . Let (f_1, \dots, f_n) be the dual basis of V^* . Choose the unique vectors $y_1, \dots, y_n \in V$ with $f_j = R_B(y_j)$ for $1 \leq j \leq n$. Then $\langle x_i, y_j \rangle = f_j(x_i) = \delta_{ij}$ for $1 \leq i, j \leq n$. Also, (y_1, \dots, y_n) is an ordered basis for V , being the image of an ordered basis under the semi-linear isomorphism R_B^{-1} .

Now let $T : V \rightarrow V$ be a fixed linear operator on the complex inner product space V . We show there is a unique linear operator $T' : V \rightarrow V$ (the *adjoint* of T relative to the complex inner product) satisfying

$$\langle T(x), y \rangle = \langle x, T'(y) \rangle \quad \text{for all } x, y \in V. \quad (13.7)$$

As in our earlier discussion of nondegenerate bilinear forms, the stated requirement on T is equivalent to the equality of functions $T^* \circ R_B = R_B \circ T'$ (note both sides are functions from V to V^*). Since R_B is a bijection, the only possible function $T' : V \rightarrow V$ that might work is $T' = R_B^{-1} \circ T^* \circ R_B$, proving uniqueness. This function is linear, being the composition of one linear map and two semi-linear maps, proving existence.

Suppose $B = (x_1, \dots, x_n)$ is an ordered basis of V , and let A be the matrix of the linear map T relative to this basis. Then for $1 \leq j \leq n$, $T(x_j) = \sum_{i=1}^n A(i, j)x_i$. Let $B^d = (y_1, \dots, y_n)$ be the dual basis for B ; we then claim that *the matrix of the linear map T' relative to this dual basis B^d is the conjugate-transpose A^** . (Recall that $A^*(i, j) = \overline{A(j, i)}$ for $1 \leq i, j \leq n$.) To check the claim, fix j and write $T'(y_j) = \sum_{k=1}^n c_{k,j}y_k$ for unique complex scalars $c_{k,j}$. Now, fix i with $1 \leq i \leq n$. On one hand,

$$\langle x_i, T'(y_j) \rangle = \left\langle x_i, \sum_{k=1}^n c_{k,j}y_k \right\rangle = \sum_{k=1}^n \overline{c_{k,j}} \langle x_i, y_k \rangle = \overline{c_{i,j}}$$

since B and B^d are dual bases. On the other hand, using (13.7) gives

$$\langle x_i, T'(y_j) \rangle = \langle T(x_i), y_j \rangle = \left\langle \sum_{k=1}^n A(k, i)x_k, y_j \right\rangle = \sum_{k=1}^n A(k, i) \langle x_k, y_j \rangle = A(j, i).$$

So the matrix of T' has i, j -entry $c_{i,j} = \overline{A(j, i)} = A^*(i, j)$, as claimed.

13.11 Comments on Infinite-Dimensional Spaces

Although many of the results in this chapter apply to arbitrary vector spaces V , some theorems require the hypothesis that V be finite-dimensional. For example, we used finite-dimensionality in the proofs that $\dim(V) = \dim(V^*)$, that $\text{ev} : V \rightarrow V^{**}$ is an isomorphism, and that \mathcal{Z} and \mathcal{A} are bijections between the subspace lattices of V and V^* . If we start with an infinite-dimensional space V , it can be shown that the dimension of the dual space V^* is a larger infinite cardinal than $\dim(V)$ (cf. Exercises 10 and 11). Then $\dim(V^{**})$ is larger still, so that V and V^{**} are never isomorphic in the infinite-dimensional case.

In order to obtain more satisfactory results for infinite-dimensional spaces, one typically adds extra structure to the vector space V and then modifies the definition of the dual space V^* to take account of this extra structure. We now sketch how this is done for Banach spaces. In many real or complex vector spaces V , it is possible to define the *norm* or *length* of a vector, denoted $\|x\|$, satisfying these conditions: $\|0\| = 0$; for all nonzero $x \in V$, $\|x\| \in \mathbb{R}^+$; $\|cx\| = |c| \cdot \|x\|$ for all $c \in F$ and $x \in V$; and $\|x + y\| \leq \|x\| + \|y\|$ for all $x, y \in V$. The space V together with the norm function on V is called a *normed vector space*. Using the

norm, one can define the *distance* between two vectors by setting $d(x, y) = \|x - y\|$ for $x, y \in V$, and this turns V into a *metric space*. If this metric space is *complete* (meaning every Cauchy sequence of vectors in V converges to a vector in V ; see Chapter 14 for more details), then V is called a *Banach space*.

Given two Banach spaces V and W , one could study all linear maps from V to W . However, it is more fruitful to restrict attention to the *continuous* linear maps from V to W . It can be shown that a linear map $T : V \rightarrow W$ is continuous iff T is continuous at 0_V iff there is a real constant $C > 0$ such that $\|T(x)\| \leq C\|x\|$ for all $x \in V$. For this reason, continuous linear maps are also called *bounded* linear operators. In this setting, we adjust the definition of the dual space so that V^* consists of all *continuous* linear maps from V to F (where F is \mathbb{R} or \mathbb{C}).

With this adjusted definition, various Banach spaces that arise in integration theory turn out to be duals of each other. For instance, for each fixed p with $1 < p < \infty$, there is a real Banach space L_p consisting of Lebesgue-measurable functions $f : \mathbb{R} \rightarrow \mathbb{R}$ such that $\int_{\mathbb{R}} |f(x)|^p dx < \infty$. Choose q so that $1/p + 1/q = 1$. In advanced analysis, it is shown that L_q is isomorphic to L_p^* as a Banach space. The isomorphism sends $g \in L_q$ to the continuous linear functional that maps $f \in L_p$ to $\int_{\mathbb{R}} f(x)g(x) dx$. It follows from this that $L_p^{**} \cong L_q^* \cong L_p$, as we might have expected from the finite-dimensional case. However, not all Banach spaces are isomorphic to their double duals. A Banach space V that is isomorphic to V^{**} is called *reflexive*.

Another subtlety that arises in Banach spaces involves the \mathcal{Z} and \mathcal{A} operators. Recall that for $S \subseteq V^*$, $\mathcal{Z}(S)$ consists of all $x \in V$ such that $f(x) = 0$ for all $f \in S$. We can express this in symbols by writing $\mathcal{Z}(S) = \bigcap_{f \in S} f^{-1}[\{0\}]$. In the current setting, each $f \in S$ must be *continuous*, so that the inverse image $f^{-1}[\{0\}]$ of the closed set $\{0\}$ is a closed subset of the metric space V . Then $\mathcal{Z}(S)$, being an intersection of closed sets, is always a closed set (as well as a subspace) in V . Similarly, one can show that V^* is a Banach space and $\mathcal{A}(T)$ is closed in V^* for all $T \subseteq V$. Thus, in order to turn the operators \mathcal{Z} and \mathcal{A} into bijective correspondences, it becomes necessary to restrict attention to the lattices of *closed* subspaces of V and V^* .

As the above discussion suggests, one often needs to add topological ingredients to obtain an effective theory for infinite-dimensional vector spaces. In particular, restricting attention to *continuous* linear maps and *closed* linear subspaces is often essential. We shall see how this works in more detail in Chapter 14, which covers metric spaces and the basic elements of Hilbert space theory.

13.12 Affine Algebraic Geometry

We end this chapter with a brief introduction to another correspondence between spaces and functions that arises in elementary algebraic geometry. We start with an algebraically closed field F such as the field of complex numbers. Let $F^n = \{(c_1, \dots, c_n) : c_i \in F\}$ be the space of all n -tuples of elements of F . In this setting, F^n is called *affine n -space over F* . Next, let R be the polynomial ring $F[x_1, \dots, x_n]$. Each formal polynomial $p \in R$ determines a polynomial function (also denoted p) from F^n to F . For instance, when $n = 3$ and $p = x_1^3 + 2x_2x_3^2$, the polynomial function $p : F^3 \rightarrow F$ is given by $p(c_1, c_2, c_3) = c_1^3 + 2c_2c_3^2$ for $c_1, c_2, c_3 \in F$.

As in the case of linear functionals, we can define a *zero-set operator* \mathcal{Z} that maps subsets

of R to subsets of F^n . Specifically, for all $S \subseteq R$, let

$$\mathcal{Z}(S) = \{(c_1, \dots, c_n) \in F^n : p(c_1, \dots, c_n) = 0 \text{ for all } p \in S\}.$$

We can regard $\mathcal{Z}(S)$ as the solution set of the (possibly infinite) system of polynomial equations $p(\vec{c}) = 0$ for $p \in S$. Any subset of F^n that has the form $\mathcal{Z}(S)$ for some $S \subseteq R$ is called an *affine variety* in F^n . Many texts write $\mathcal{V}(S)$ for $\mathcal{Z}(S)$ and call \mathcal{Z} the *variety operator*.

We list without proof some properties of \mathcal{Z} . First, \mathcal{Z} reverses inclusions: if $S \subseteq T \subseteq R$, then $\mathcal{Z}(T) \subseteq \mathcal{Z}(S)$. Second, suppose $S \subseteq R$ generates an ideal

$$I = \{r_1 s_1 + \dots + r_m s_m : m \in \mathbb{N}^+, r_i \in R, s_i \in S\},$$

and define $\sqrt{I} = \{p \in R : p^k \in I \text{ for some } k \in \mathbb{N}^+\}$. Then $\mathcal{Z}(S) = \mathcal{Z}(I) = \mathcal{Z}(\sqrt{I})$. Third, $\mathcal{Z}(\{0\}) = F^n$ and $\mathcal{Z}(\{1\}) = \mathcal{Z}(R) = \emptyset$. Fourth, for a family of subsets $\{S_j : j \in J\}$ of R , $\mathcal{Z}(\bigcup_{j \in J} S_j) = \bigcap_{j \in J} \mathcal{Z}(S_j)$. Fifth, for a family of ideals $\{I_j : j \in J\}$ of R , $\mathcal{Z}(\sum_{j \in J} I_j) = \bigcap_{j \in J} \mathcal{Z}(I_j)$. Sixth, for ideals I_1, I_2, \dots, I_k in R , $\mathcal{Z}(I_1 I_2 \cdots I_k) = \mathcal{Z}(\bigcap_{j \in J} I_j) = \bigcup_{j \in J} \mathcal{Z}(I_j)$. The *Hilbert basis theorem* asserts that any ideal I in $R = F[x_1, \dots, x_n]$ can be generated by a finite set S . Combining this theorem with the second property, we see that any affine variety can be defined as the solution set of a *finite* system of polynomial equations.

Next we define an analogue of the annihilator operator \mathcal{A} , which maps subsets of F^n to subsets of R . Given $S \subseteq F^n$, let

$$\mathcal{A}(S) = \{p \in R : p(c_1, \dots, c_n) = 0 \text{ for all } (c_1, \dots, c_n) \in S\}.$$

It can be checked that $\mathcal{A}(S)$ is always a *radical ideal* of R ; this is an ideal I such that for all $p \in R$ and all $k \in \mathbb{N}^+$, if $p^k \in I$ then $p \in I$. Some texts write $\mathcal{I}(S)$ for $\mathcal{A}(S)$.

We state some properties of the annihilator operator and its relation to the zero-set operator. First, $\mathcal{A}(\emptyset) = R$ and $\mathcal{A}(F^n) = \{0_R\}$. Second, \mathcal{A} reverses inclusions: if $S \subseteq T \subseteq F^n$, then $\mathcal{A}(T) \subseteq \mathcal{A}(S)$. Third, for any family of subsets $\{S_j : j \in J\}$ of F^n , $\mathcal{A}(\bigcup_{j \in J} S_j) = \bigcap_{j \in J} \mathcal{A}(S_j)$. Fourth, for any $S \subseteq F^n$, $\mathcal{Z}(\mathcal{A}(S))$ is the intersection of all varieties in F^n that contain S , which is the smallest variety containing S . Fifth, for any subset S of R , $\mathcal{A}(\mathcal{Z}(S)) = \sqrt{I}$, where I is the ideal generated by S . (The last property, which requires algebraic closure of the field F , is one version of a difficult theorem of algebraic geometry called the *Nullstellensatz*.)

It follows from the results stated above that \mathcal{Z} and \mathcal{A} restrict to give poset anti-isomorphisms between the lattice of affine varieties in F^n and the lattice of radical ideals of $R = F[x_1, \dots, x_n]$. This “ideal-variety correspondence” has many additional structural properties: for instance, maximal ideals of R correspond to individual points in F^n , prime ideals of R correspond to irreducible varieties in F^n , and so on. For more details, we refer the reader to the excellent text by Cox, Little, and O’Shea [11].

13.13 Summary

Let F be a field, and let V and W be finite-dimensional vector spaces over F . Table 13.1 reviews some of the main definitions introduced in this chapter.

1. *Universal Mapping Property and Dual Bases.* For every ordered basis $B = (x_1, \dots, x_n)$ of V and every list (y_1, \dots, y_n) of vectors in W , there exists a unique F -linear map $T : V \rightarrow W$ with $T(x_i) = y_i$ for $1 \leq i \leq n$. It follows that

TABLE 13.1

Definitions of Concepts Related to Dual Spaces.

Concept	Definition
$\text{Hom}_F(V, W)$	set of all F -linear maps $T : V \rightarrow W$ (a vector space under pointwise operations)
dual space V^*	$V^* = \text{Hom}_F(V, F)$
dual basis of ordered basis $B = (x_1, \dots, x_n)$ of V	unique ordered basis $B^d = (f_1, \dots, f_n)$ of V^* with $f_i(x_j) = 1$ for $i = j$, $f_i(x_j) = 0$ for $i \neq j$
zero-set $\mathcal{Z}(S)$ of $S \subseteq V^*$	subspace $\{x \in V : f(x) = 0 \text{ for all } f \in S\}$
annihilator $\mathcal{A}(T)$ of $T \subseteq V$	subspace $\{g \in V^* : g(y) = 0 \text{ for all } y \in T\}$
double dual V^{**}	$V^{**} = \text{Hom}_F(V^*, F) = \text{Hom}_F(\text{Hom}_F(V, F), F)$
evaluation map E_x for $x \in V$	$E_x \in V^{**}$ sends $g \in V^*$ to $g(x) \in F$
$\text{ev} : V \rightarrow V^{**}$	$\text{ev}(x) = E_x$ for $x \in V$ (injective linear map)
dual map of $T : V \rightarrow W$	$T^* : W^* \rightarrow V^*$ given by $T^*(g) = g \circ T$ for $g \in W^*$
bilinear form B on V	map $B : V \times V \rightarrow F$ linear in each separate variable
$L_B : V \rightarrow V^*$ (B bilinear)	$L_B(x)$ sends y to $B(x, y)$ for $x, y \in V$
$R_B : V \rightarrow V^*$ (B bilinear)	$R_B(y)$ sends x to $B(x, y)$ for $x, y \in V$
symmetric bilinear form	$B(x, y) = B(y, x)$ for all $x, y \in V$
nondegenerate bilinear form	L_B and R_B are isomorphisms $V \cong V^*$
S^\perp for $S \subseteq V$ (relative to B)	subspace $\{y \in V : B(x, y) = 0 \text{ for all } x \in S\}$
real inner product	symm. bilinear form on real vector space V with $\langle x, x \rangle > 0$ for all $x \neq 0$ in V
adjoint map of $T : V \rightarrow V$ (relative to B)	unique linear map $T' : V \rightarrow V$ with $B(T(x), y) = B(x, T'(y))$ for $x, y \in V$
semi-linear map $T : V \rightarrow W$ ($F = \mathbb{C}$)	$T(x + y) = T(x) + T(y)$ for all $x, y \in V$, and $T(cx) = \bar{c}T(x)$ for all $x \in V, c \in \mathbb{C}$

there exists a unique dual basis $B^d = (f_1, \dots, f_n)$ of V^* with $f_i(x_j) = \delta_{i,j}$ for all $i, j \in [n]$. Conversely, for any ordered basis $C = (g_1, \dots, g_n)$ of V^* , there is a unique ordered basis $C^D = (z_1, \dots, z_n)$ of V with $g_i(z_j) = \delta_{i,j}$ for all $i, j \in [n]$. For finite-dimensional V , $\dim(V) = \dim(V^*)$ and $V \cong V^*$, but in general there is no *natural* isomorphism between V and V^* .

2. *Zero-Set Operator and Annihilator Operator.* For $S \subseteq V^*$, $\mathcal{Z}(S)$ is the set of $x \in V$ such that $f(x) = 0$ for all $f \in S$. For $T \subseteq V$, $\mathcal{A}(T)$ is the set of $g \in V^*$ such that $g(y) = 0$ for all $y \in T$. The operators \mathcal{Z} and \mathcal{A} reverse inclusions, map subsets to subspaces, and satisfy $\mathcal{Z}(S) = \mathcal{Z}(\langle S \rangle)$, $\mathcal{A}(T) = \mathcal{A}(\langle T \rangle)$, $\mathcal{A}(\mathcal{Z}(S)) = \langle S \rangle$, and $\mathcal{Z}(\mathcal{A}(T)) = \langle T \rangle$. So, for finite-dimensional V , \mathcal{Z} and \mathcal{A} restrict to give poset anti-isomorphisms between the lattice of subspaces of V and the lattice of subspaces of V^* . For subspaces U and W of V^* , $\dim(U) + \dim(\mathcal{Z}(U)) = \dim(V)$, $\mathcal{Z}(U + W) = \mathcal{Z}(U) \cap \mathcal{Z}(W)$, $\mathcal{Z}(U \cap W) = \mathcal{Z}(U) + \mathcal{Z}(W)$, and similar formulas hold for \mathcal{A} .
3. *Double Duals.* There is an injective linear map $\text{ev} : V \rightarrow V^{**}$ sending each $x \in V$ to “evaluation at x ,” which is the map $E_x : V^* \rightarrow F$ given by $E_x(g) = g(x)$ for $g \in V^*$. For finite-dimensional V , ev gives a *natural* vector-space isomorphism $V \cong V^{**}$, meaning that $T^{**} = \text{ev}_W \circ T \circ \text{ev}_V^{-1}$ for all linear $T : V \rightarrow W$. Using ev , the zero-set and annihilator operators linking V^* and V^{**} can be identified with the annihilator and zero-set operators (respectively) linking V and V^* .
4. *Dual Maps.* For each linear map $T : V \rightarrow W$, there is a linear dual map

$T^* : W^* \rightarrow V^*$ given by $T^*(g) = g \circ T$ for $g \in W^*$. If A is the matrix of T relative to the ordered bases B for V and C for W , then the matrix of T^* relative to the dual bases C^d and B^d is the transpose of A . For linear maps T and U and $c \in F$, the formulas $(T + U)^* = T^* + U^*$, $(cT)^* = c(T^*)$, $(\text{id}_V)^* = \text{id}_{V^*}$, and $(T \circ U)^* = U^* \circ T^*$ hold when defined.

5. *Nondegenerate Symmetric Bilinear Forms.* Let B be a nondegenerate symmetric bilinear form. Since $L_B = R_B : V \rightarrow V^*$ is an isomorphism, every linear functional $f : V \rightarrow F$ has the form $f(x) = B(x, y) = B(y, x)$ for a unique $y \in V$. Every linear map $T : V \rightarrow V$ has a unique adjoint $T' : V \rightarrow V$ satisfying $B(T(x), y) = B(x, T'(y))$ for all $x, y \in V$; in fact, $T' = R_B^{-1} \circ T^* \circ R_B$. For every subset S of V , $S^\perp = \{z \in V : B(x, z) = 0 \text{ for all } x \in S\}$ is a subspace of V . The map $W \mapsto W^\perp$ is a poset anti-isomorphism on the lattice of subspaces of V with $\dim(W) + \dim(W^\perp) = \dim(V)$, $W^{\perp\perp} = W$, $(W \cap Z)^\perp = W^\perp + Z^\perp$, and $(W + Z)^\perp = W^\perp \cap Z^\perp$. However, it is not always the case that $W \cap W^\perp = \{0\}$ and $V = W \oplus W^\perp$.
6. *Real Inner Product Spaces.* For a real vector space V , an inner product is a symmetric bilinear form on V satisfying $\langle x, x \rangle > 0$ for all $x \neq 0$, which guarantees nondegeneracy of the form. All results of the preceding item apply, and now it is true that $W \cap W^\perp = \{0\}$ and $V = W \oplus W^\perp$ for all subspaces W of V .
7. *Complex Inner Product Spaces.* For a complex vector space V , a complex inner product is a map from $V \times V$ to \mathbb{C} that is linear in the first argument, has the conjugate-symmetry $\langle x, y \rangle = \overline{\langle y, x \rangle}$ for $x, y \in V$, and satisfies $\langle x, x \rangle \in \mathbb{R}^+$ for all nonzero $x \in V$. Results similar to the preceding two items hold, but now $L_B : V \rightarrow V^{*s}$ is an isomorphism between V and the space V^{*s} of semi-linear maps from V to \mathbb{C} , and $R_B : V \rightarrow V^*$ is a semi-isomorphism. Properties of the map $W \mapsto W^\perp$ are the same as for real inner product spaces.
8. *Infinite-Dimensional Spaces.* For infinite-dimensional F -vector spaces with topological structure (e.g., Banach spaces where each vector has a norm), one redefines V^* to consist of *continuous* linear maps from V to F . Here, V^{**} may or may not be isomorphic to V , and one must restrict to *closed* linear subspaces to make \mathcal{Z} and \mathcal{A} be bijective.

13.14 Exercises

Unless otherwise stated, assume F is a field and V, W are finite-dimensional F -vector spaces in these exercises.

1. Let S be a set, and let Z be the set of functions $f : S \rightarrow W$. (a) Verify that Z is a commutative group under pointwise addition of functions. (b) Verify that Z is an F -vector space under pointwise operations.
2. (a) Let V and W be F -vector spaces. Prove carefully that $\text{Hom}_F(V, W)$ is a subspace of the vector space of all functions from V to W (under pointwise operations). (b) Explain why the set V of all differentiable functions $f : [0, 1] \rightarrow \mathbb{R}$ is a subspace of the set of all functions from $[0, 1]$ to \mathbb{R} (under pointwise operations).
3. Let V be the vector space of differentiable functions $f : [0, 1] \rightarrow \mathbb{R}$ under pointwise

operations. Decide whether each operator below is a linear functional in V^* .

- (a) $D : V \rightarrow V$ given by $D(f) = f'$; (b) $E : V \rightarrow \mathbb{R}$ given by $E(f) = f(1/3)$;
- (c) $F : V \rightarrow \mathbb{R}$ given by $F(f) = f'(1/2)$; (d) $G : V \rightarrow \mathbb{R}$ given by $G(f) = f(0)f(1)$;
- (e) $H : V \rightarrow \mathbb{R}$ given by $H(f) = \int_0^1 f(x) dx$; (f) $I : V \rightarrow \mathbb{R}$ given by $I(f) = \int_0^1 xf(x^2) dx$.

4. Fix an integer $n > 1$. Which of the following functions with domain $M_n(F)$ are in $M_n(F)^*$? Explain. (a) the map sending A to $A(1, 2)$; (b) the trace map given by $\text{tr}(A) = \sum_{i=1}^n A(i, i)$; (c) the determinant function sending A to $\det(A)$; (d) the transpose map sending A to A^T ; (e) the map sending A to the sum of all entries of A ; (f) for fixed $v \in M_{n,1}(F)$, the map sending A to $v^T Av$; (g) for $F = \mathbb{R}$, the map sending A to the number of zero entries of A .
5. Check that the map T defined in the first paragraph of §13.2 is F -linear and sends each x_i to y_i .
6. Fix $m, n \in \mathbb{N}^+$. Let $B = \{E_{i,j} : 1 \leq i \leq m, 1 \leq j \leq n\}$ be the standard basis for $M_{m,n}(F)$, where $E_{i,j}$ is the matrix with 1 in the i, j -position and zeroes elsewhere. Explicitly describe how each element in the dual basis B^d acts on an arbitrary matrix $A \in M_{m,n}(F)$.
7. (a) For $n \in \mathbb{N}^+$, define a map $T : M_{1,n}(F) \rightarrow M_{n,1}(F)^*$ by letting $T(w)$ (for a given row vector w) be the map $v \mapsto wv$ (for all column vectors v). Prove that T is a vector space isomorphism. (We can use this isomorphism to identify the dual space of the vector space of $n \times 1$ column vectors with the vector space of $1 \times n$ row vectors.) (b) Let $(E(i, 1) : 1 \leq i \leq n)$ be the standard ordered basis of $M_{n,1}(F)$, and let $(G(j) : 1 \leq j \leq n)$ be the dual basis of $M_{n,1}(F)^*$. Compute $T^{-1}(G(j))$ for each j .
8. (a) Find a dual basis for the ordered basis

$$B = \left(\begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 3 \end{bmatrix}, \begin{bmatrix} -3 \\ 4 \\ -1 \end{bmatrix} \right)$$

of \mathbb{R}^3 , using Exercise 7 to describe the final answer as a list of row vectors. (b) Given an ordered basis (v_1, \dots, v_n) of F^n (viewed as column vectors), describe how to use matrix inversion to compute the dual basis of $(F^n)^*$ (viewed as row vectors).

9. This exercise shows what can go wrong with the dual basis construction in the case of infinite-dimensional vector spaces. Let V be the vector space $F[x]$ of polynomials with coefficients in F . (a) For each $k \in \mathbb{N}$, let $C_k : F[x] \rightarrow F$ be the map given by $C_k(\sum_{i \geq 0} a_i x^i) = a_k$ for $a_i \in F$. Prove each $C_k \in F[x]^*$. (b) We know $\{x^j : j \in \mathbb{N}\}$ is a basis of $F[x]$. Compute $C_k(x^j)$ for all $j, k \in \mathbb{N}$. (c) Prove $\{C_k : k \in \mathbb{N}\}$ is F -linearly independent. (d) Define $h(\sum_{i=0}^n a_i x^i) = \sum_{i=0}^n a_i$. Prove $h \in F[x]^*$, but h is outside the span of $\{C_k : k \in \mathbb{N}\}$. (e) Define $C_k : F[[x]] \rightarrow F$ by the same formula as in (a). Is $\{C_k : k \in \mathbb{N}\}$ a basis of $F[[x]]^*$?
10. (a) Formulate and prove a generalization of the universal mapping property in §13.2 that applies to an infinite-dimensional vector space V with infinite basis B . (b) Use the UMP to prove that the dual space V^* is isomorphic to the vector space of all functions from B to F with pointwise operations.
11. Prove that $F[x]^*$ and $F[[x]]$ are isomorphic as F -vector spaces.

12. *Elementary Operations and Dual Bases.* Let $B = (x_1, \dots, x_n)$ be an ordered basis of an F -vector space V with dual basis $B^d = (g_1, \dots, g_n)$. (a) Suppose C is the ordered basis obtained from B by switching the positions of x_i and x_j . Describe (with proof) how C^d is related to B^d . (b) Suppose C is obtained from B by multiplying x_i by $d \neq 0$ in F . How is C^d related to B^d ? (c) Suppose C is obtained from B by replacing x_i by $x_i + bx_j$ for some $i \neq j$ and some $b \in F$. How is C^d related to B^d ?
13. Let Re and Im be the functions from \mathbb{C} to \mathbb{R} given by $\text{Re}(a + ib) = a$ and $\text{Im}(a + ib) = b$ for $a, b \in \mathbb{R}$. Find an ordered basis $B = (z_1, z_2)$ of \mathbb{C} (viewed as a real vector space) such that $B^d = (\text{Im}, \text{Re})$.
14. (a) Define $f, g \in (\mathbb{R}^2)^*$ by $f((x, y)) = x + y$ and $g((x, y)) = 2x - 3y$ for $x, y \in \mathbb{R}$. Find an ordered basis $B = (v_1, v_2)$ of \mathbb{R}^2 such that $B^d = (f, g)$. (b) Given an ordered basis $C = (f_1, \dots, f_n)$ of $(F^n)^*$ (viewed as row vectors), describe a matrix computation that will produce an ordered basis $B = (x_1, \dots, x_n)$ of F^n (viewed as column vectors) such that $C = B^d$.
15. Define $f, g, h \in (\mathbb{R}^4)^*$ by setting $f((x_1, x_2, x_3, x_4)) = x_1 + 2x_2 - 3x_3 + x_4$, $g((x_1, x_2, x_3, x_4)) = 4x_1 - x_2 - x_4$, and $h((x_1, x_2, x_3, x_4)) = 3x_2 + 5x_3 - x_4$ for $x_i \in \mathbb{R}$. (a) Find a basis for $\mathcal{Z}(\{f, g, h\})$. (b) Find a basis for $\mathcal{Z}(\{f, g\})$. (c) Find a basis for $\mathcal{Z}(\{f\})$.
16. Let V be the vector space of differentiable functions $f : [0, 1] \rightarrow \mathbb{R}$. Define $D, I \in V^*$ by $D(f) = f'(1/2)$ and $I(f) = \int_0^1 f(x) dx$ for $f \in V$. Describe geometrically which functions f belong to $\mathcal{Z}(\{D, I\})$.
17. Let $B = (x_1, \dots, x_n)$ be an ordered basis of V with dual basis $B^d = (f_1, \dots, f_n)$. (a) Prove: for all $v \in V$, $v = \sum_{i=1}^n f_i(v)x_i$. (b) For $1 \leq i \leq n$, explicitly describe $\mathcal{Z}(\{f_1, \dots, f_i\})$ and $\mathcal{A}(\{x_1, \dots, x_i\})$.
18. (a) Negate the definition to give a useful completion of this sentence: for all $g \in V^*$, $g \notin \mathcal{A}(T)$ iff ... (b) Carefully prove that $\mathcal{A}(\emptyset) = V^*$.
19. For each subset T of \mathbb{R}^n , describe $\mathcal{A}(T)$ in $(\mathbb{R}^n)^*$ as explicitly as possible. (a) $n = 2$, $T = \{(2, 3)\}$; (b) $n = 3$, $T = \{(s, s, s) : s \in \mathbb{R}\}$; (c) $n = 4$, $T = \{(w, x, y, z) : w + x + y + z = 0$ and $w + z = x + y\}$; (d) $n = 2$, $T = \{(x, y) : x^2 + y^2 = 1\}$; (e) $n = 3$, $T = \{(x, y, z) : x^2 + y^2 = 1, z = 0\}$.
20. For each T in Exercise 19, describe $\mathcal{Z}(\mathcal{A}(T))$.
21. Let $n > 1$. For each subset T of $M_n(\mathbb{R})$, describe $\mathcal{A}(T)$ in $M_n(\mathbb{R})^*$ as explicitly as possible. (a) $T = \{A \in M_n(\mathbb{R}) : \text{tr}(A) = 0\}$; (b) $T = \{A \in M_n(\mathbb{R}) : \det(A) = 0\}$; (c) $T = \{\text{upper-triangular } A \in M_n(\mathbb{R})\}$; (d) $T = \{\text{diagonal } A \in M_n(\mathbb{R})\}$; (e) $T = \{cI_n : c \in \mathbb{R}\}$.
22. (a) Prove: for an n -dimensional vector space V and all $T \subseteq V^*$, there exists a subset S of V^* of size at most n with $\mathcal{Z}(T) = \mathcal{Z}(S)$. (b) Formulate and prove a statement analogous to (a) for annihilators.
23. (a) Which of the five properties of zero sets listed in §13.3 remain true for infinite-dimensional vector spaces V ? (b) Which of the five properties of annihilators listed in §13.4 remain true for infinite-dimensional vector spaces V ?
24. (a) Show directly from the definitions that $\mathcal{A}(\mathcal{Z}(S)) \supseteq S$ for all $S \subseteq V^*$. (b) Use (a) and the identification of V and V^{**} to deduce $\mathcal{Z}(\mathcal{A}(T)) \supseteq T$ for all $T \subseteq V$.
25. (a) Prove: for any family of subsets $\{S_i : i \in I\}$ of V^* , $\mathcal{Z}(\bigcup_{i \in I} S_i) = \bigcap_{i \in I} \mathcal{Z}(S_i)$. (b) Formulate and prove a statement similar to (a) for the \mathcal{A} operator. (c) For arbitrary subsets $S, T \subseteq V^*$, must it be true that $\mathcal{Z}(S \cap T) = \mathcal{Z}(S) + \mathcal{Z}(T)$?

26. Use the theory of dual spaces to prove that every k -dimensional subspace W of F^n is the solution set of a system of $n - k$ homogeneous linear equations. Also prove that W is the intersection of $n - k$ linear hyperplanes (subspaces of dimension $n - 1$).
27. Let $F = \mathbb{Z}_2$ and $V = F^2$. (a) Explicitly describe all elements of V^* by giving the domain, codomain, and values of each $f \in V^*$. (b) Similarly, describe all elements $E \in V^{**}$. (c) For each E in (b), find an $x \in V$ such that $E = E_x$ (evaluation at x). (d) List all ordered bases of V . For each basis, find the associated dual basis of V^* .
28. Let $v_1 = (1, 2, 3)$, $v_2 = (0, 1, -1)$, and $v_3 = (0, 0, 1)$, so $B = (v_1, v_2, v_3)$ is an ordered basis of \mathbb{R}^3 . (a) Let $B^d = (f_1, f_2, f_3)$ be the dual basis of B . Compute the matrix of each f_i relative to the standard ordered bases of \mathbb{R}^3 and \mathbb{R} . (b) Let $w = (2, 3, -1)$. Find $E_w(f_i)$ for $i = 1, 2, 3$. (c) Use the UMP to obtain an element $E \in (\mathbb{R}^3)^{**}$ that sends f_1 to 4, f_2 to -1 , and f_3 to 0. Find $x \in \mathbb{R}^3$ such that $E = E_x$.
29. Let $B = (x_1, \dots, x_n)$ be an ordered basis for V , and let $B^d = (f_1, \dots, f_n)$ be the dual basis for V^* . Show that $(E_{x_1}, \dots, E_{x_n})$ is the dual basis of B^d in V^{**} .
30. Prove that if V is infinite-dimensional, then the map $\text{ev} : V \rightarrow V^{**}$ is not surjective. (Use Exercise 10.)
31. (a) Give a direct proof (using bases but not facts about \mathcal{A}) that $\dim(W) + \dim(\mathcal{Z}(W)) = \dim(V)$ for all subspaces W of V^* . (b) Use (a) to deduce $\dim(Y) + \dim(\mathcal{A}(Y)) = \dim(V)$ for all subspaces Y of V .
32. (a) In §13.6, give the details of the proof that $\mathcal{A}(\mathcal{Z}(W)) = W$ for all $W \in \mathcal{S}(V^*)$. (b) Explain why $\mathcal{Z}(\mathcal{A}(T)) = \langle T \rangle$ for all $T \subseteq V$.
33. Let (e_1, e_2, e_3) be the standard ordered basis of \mathbb{R}^3 , and let (f_1, f_2, f_3) be the dual basis of $(\mathbb{R}^3)^*$. Let $X = \langle f_1, f_2 \rangle$ and $Y = \langle f_2, f_3 \rangle$. Verify (13.4) by explicitly computing $X \cap Y$, $X + Y$, $\mathcal{Z}(X)$, $\mathcal{Z}(Y)$, $\mathcal{Z}(X) + \mathcal{Z}(Y)$, $\mathcal{Z}(X) \cap \mathcal{Z}(Y)$, $\mathcal{Z}(X \cap Y)$, and $\mathcal{Z}(X + Y)$ in \mathbb{R}^3 .
34. Let $X = \langle (1, 2, 1) \rangle$ and $Y = \langle (2, -1, 0) \rangle$ in \mathbb{R}^3 . Verify (13.5) by explicitly computing $X \cap Y$, $X + Y$, $\mathcal{A}(X)$, $\mathcal{A}(Y)$, $\mathcal{A}(X) + \mathcal{A}(Y)$, $\mathcal{A}(X) \cap \mathcal{A}(Y)$, $\mathcal{A}(X \cap Y)$, and $\mathcal{A}(X + Y)$ in $(\mathbb{R}^3)^*$.
35. Let $T = \{(1, 0, 1), (2, 2, 2)\}$ in \mathbb{R}^3 . Find $\mathcal{Z}(\mathcal{A}(T))$: (a) by using a theorem; (b) by using the definitions of \mathcal{A} and \mathcal{Z} .
36. Let $\mathcal{S}(V)$ be the poset of subspaces of V , ordered by set inclusion. (a) Prove $\inf(X, Y) = X \cap Y$ for $X, Y \in \mathcal{S}(V)$. (b) Prove $\sup(X, Y) = X + Y$ for $X, Y \in \mathcal{S}(V)$. (c) Prove $\inf(X_i : i \in I) = \bigcap_{i \in I} X_i$ for a family $X_i \in \mathcal{S}(V)$. (d) Prove $\sup(X_i : i \in I) = \sum_{i \in I} X_i$ for $X_i \in \mathcal{S}(V)$.
37. Suppose U and V are lattices and $f : U \rightarrow V$ is a bijection such that for all $x, y \in U$, $x \leq y$ iff $f(y) \leq f(x)$. (a) Prove $f(\inf(x, y)) = \sup(f(x), f(y))$ and $f(\sup(x, y)) = \inf(f(x), f(y))$ for all $x, y \in U$. (b) Explain how (13.4) and (13.5) follow from (a). (c) Try to prove (13.4) directly from the definitions, without explicitly invoking (a).
38. State and prove analogues of (13.4) and (13.5) that involve sup's and inf's of an indexed collection of subspaces.
39. Let F be a finite field. Show that, for $0 \leq k \leq n$, the number of k -dimensional subspaces of an n -dimensional F -vector space V is the same as the number of $(n - k)$ -dimensional subspaces of V .

40. Assume $\dim(V) = n > 1$. Prove there does not exist any isomorphism $f : V \rightarrow V^*$ such that $f \circ T = T^* \circ f$ for all linear maps $T : V \rightarrow V$.
41. (a) Prove $(V \times W)^* \cong V^* \times W^*$. (b) Under the isomorphism in (a), prove $\mathcal{A}(V \times \{0\})$ maps to $\{0\} \times W^*$. (c) Given $k \in \mathbb{N}^+$ and vector spaces V_1, \dots, V_k , prove that $(V_1 \times V_2 \times \dots \times V_k)^*$ is isomorphic to $V_1^* \times V_2^* \times \dots \times V_k^*$.
42. Suppose a linear map $T : V \rightarrow W$ has matrix A relative to an ordered basis B of V and an ordered basis C of W . What is the matrix of T^{**} relative to the bases B^{dd} and C^{dd} ?
43. Let $D : F[x] \rightarrow F[x]$ be the map given by $D(\sum_{k \geq 0} a_k x^k) = \sum_{k \geq 1} k a_k x^{k-1}$ for $a_k \in F$. Describe the action of D^* on the vector space $F[x]^* \cong F[[x]]$ (see Exercise 11).
44. (a) Prove: for all linear $T : V \rightarrow V$ and all $k \in \mathbb{N}$, $(T^k)^* = (T^*)^k$. (b) Prove the formula in (a) holds for all $k \in \mathbb{Z}$ if T is invertible. (c) Suppose two linear maps $T, U : V \rightarrow V$ are similar. Prove T^* and U^* are similar.
45. Suppose $T : V \rightarrow W$ is a linear map. (a) Prove: for all $S \subseteq W$, $T^*[\mathcal{A}(S)] \subseteq \mathcal{A}(T^{-1}[S])$. (b) Give an example where equality does not hold in (a). (c) Prove: for all $S \subseteq V$, $\mathcal{A}(T[S]) = (T^*)^{-1}[\mathcal{A}(S)]$.
46. Let $T : V \rightarrow W$ be a linear map. (a) Prove $\ker(T^*) = \mathcal{A}(\text{img}(T))$. (b) Deduce that T^* is injective iff T is surjective.
47. Let $T : V \rightarrow W$ be a linear map. (a) Prove $\text{img}(T^*) = \mathcal{A}(\ker(T))$. (b) Deduce that T^* is surjective iff T is injective.
48. Prove results similar to Exercises 46 and 47 for the adjoint of a linear operator on a real or complex inner product space.
49. Let B be a bilinear form on V . (a) Verify that $L_B(cx) = cL_B(x)$ for $x \in V$ and $c \in F$. (b) Verify that $R_B : V \rightarrow V^*$ maps into V^* and is linear.
50. Let $A \in M_n(F)$ be a fixed matrix. (a) Show that $B : F^n \times F^n \rightarrow F$, given by $B(v, w) = v^T A w$ for $v, w \in F^n$, is a bilinear form. (b) Show that B is a symmetric form iff $A^T = A$. (c) Show that B is nondegenerate on the left iff B is nondegenerate on the right iff $\det(A) \neq 0$.
51. *UMP for Bilinear Forms.* Let (x_1, \dots, x_n) be an ordered basis for V . (a) Prove that for every $A \in M_n(F)$, there exists a unique bilinear form B_A on V such that $B_A(x_i, x_j) = A(i, j)$ for $1 \leq i, j \leq n$, namely

$$B_A \left(\sum_{i=1}^n c_i x_i, \sum_{j=1}^n d_j x_j \right) = \sum_{i=1}^n \sum_{j=1}^n c_i A(i, j) d_j.$$

- (b) Conversely, given any bilinear form B on V , define $A \in M_n(F)$ by $A(i, j) = B(x_i, x_j)$ for $1 \leq i, j \leq n$. Show that $B = B_A$. (c) Show that B_A is nondegenerate on the left iff B_A is nondegenerate on the right iff $\det(A) \neq 0$. (d) Show that B_A is a symmetric form iff A is a symmetric matrix.
52. (a) Give an example of a degenerate symmetric bilinear form B on \mathbb{R}^3 with $B(e_i, e_j) \neq 0$ for all $i, j \in \{1, 2, 3\}$. Find all $v \in \mathbb{R}^3$ with $L_B(v) = 0$. (b) Give an example of a nondegenerate bilinear form B on \mathbb{R}^4 with $B(x, x) = 0$ for all $x \in \mathbb{R}^4$. (c) Suppose B is a symmetric bilinear form on a real vector space V such that $B(v, v) = 0$ for all $v \in V$. Prove $B = 0$.

53. Consider the following bilinear forms on \mathbb{R}^3 :

$$\begin{aligned}B_1(x, y) &= x_1y_1 + x_2y_2 + x_3y_3; \\B_2(x, y) &= x_1y_2 - x_2y_1 + 2x_3y_3; \\B_3(x, y) &= x_1y_1 + x_1y_2 + x_1y_3 + x_2y_2 + x_2y_3 + x_3y_3; \\B_4(x, y) &= x_1y_1 + x_1y_3 + 2x_2y_2 + x_3y_1 + x_3y_3.\end{aligned}$$

- (a) Which of these bilinear forms are symmetric? Which are nondegenerate?
 (b) Let $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ be the linear functional $f(x) = 3x_1 - 2x_2 - x_3$ for $x \in \mathbb{R}^3$. For $1 \leq i \leq 4$, find $y \in \mathbb{R}^3$ such that $f = R_{B_i}(y)$ or explain why this cannot be done.
 (c) For f given in (b) and $1 \leq i \leq 4$, find $z \in \mathbb{R}^3$ such that $f = L_{B_i}(z)$ or explain why this cannot be done.
54. Let V have ordered basis $X = (x_1, \dots, x_n)$, let the dual basis of V^* be $X^d = (f_1, \dots, f_n)$, let $A \in M_n(F)$, and let B_A be the bilinear form on V defined in Exercise 51. (a) What is the matrix of the linear map $R_B : V \rightarrow V^*$ relative to the basis X of V and the basis X^d of V^* ? (b) What is the matrix of the linear map $L_B : V \rightarrow V^*$ relative to the basis X of V and the basis X^d of V^* ? (c) Assume B_A is symmetric and nondegenerate, and $T : V \rightarrow V$ is a linear map with matrix C relative to X . What is the matrix of the adjoint map T' relative to X ?
55. Let $T : \mathbb{R}^4 \rightarrow \mathbb{R}^4$ be the map $T(x) = Ax$ for $x \in \mathbb{R}^4$, where
- $$A = \begin{bmatrix} 2 & 0 & -1 & 3 \\ 1 & 1 & -2 & -1 \\ 0 & 1 & 1 & 2 \\ 3 & 5 & -6 & 2 \end{bmatrix}.$$
- (a) Find the adjoint of T relative to the standard inner product on \mathbb{R}^4 . (b) Find the adjoint of T relative to the symmetric bilinear form $B(x, y) = x_1y_1 + x_2y_2 + x_3y_3 - x_4y_4$ for $x, y \in \mathbb{R}^4$. (c) Find the adjoint of T relative to the symmetric bilinear form $B(x, y) = x_1y_2 + x_2y_1 + 2x_3y_4 + 2x_4y_3$ for $x, y \in \mathbb{R}^4$.
56. Given vector spaces V and W , a function $B : V \times W \rightarrow F$ is *bilinear* iff for all $v \in V$ and $w \in W$, the functions $B(v, \cdot)$ and $B(\cdot, w)$ (obtained by fixing one of the variables) are F -linear. (a) Use B to define F -linear maps $L_B : V \rightarrow W^*$ and $R_B : W \rightarrow V^*$. (b) Define $B : V \times V^* \rightarrow F$ by $B(v, f) = f(v)$ for $v \in V$ and $f \in V^*$. Prove B is bilinear. For this choice of B , what are the maps L_B and R_B defined in (a)?
57. Let $U = \{(t, 3t, t) : t \in \mathbb{R}\}$ and $W = \{(x, y, z) \in \mathbb{R}^3 : 5x - 2y + 3z = 0\}$. Compute U^\perp and W^\perp relative to the standard inner product on \mathbb{R}^3 .
58. For each linear operator T on V , let T' denote the adjoint operator relative to a fixed nondegenerate symmetric bilinear form B . (a) Use the defining formula (13.6) to prove: for all $S, T \in \text{Hom}_F(V, V)$ and $c \in F$, $(S+T)' = S' + T'$, $(cS)' = c(S')$, and $(S \circ T)' = T' \circ S'$. (b) Reprove the formulas in (a) using the relation between T' and T^* .
59. Give an example of a vector space V , a symmetric nondegenerate bilinear form B on V , and a subspace W of V such that $W \cap W^\perp \neq \{0\}$ and $V \neq W \oplus W^\perp$. For your choice of W , verify that $W^{\perp\perp} = W$ and $\dim(W) + \dim(W^\perp) = \dim(V)$.
60. Let $T : V \rightarrow W$ and $U : W \rightarrow Z$ be semi-linear maps between complex vector spaces V , W , and Z . Prove: (a) T is injective iff $\ker(T) = \{0\}$. (b) $U \circ T$ is

\mathbb{C} -linear. (c) The composition of T with a linear map (in either order) is semi-linear. (d) If T is bijective, then T maps linearly independent sets to linearly independent sets and bases to bases. (e) If T is bijective, then T^{-1} is a semi-isomorphism. (f) The image of a subspace under T or T^{-1} is another subspace, which has the same dimension as the original subspace if T is a semi-isomorphism. (g) The rank-nullity formula holds for T .

61. Prove that for a finite-dimensional complex vector space V , V^{*s} is a complex vector space with $\dim(V^{*s}) = \dim(V)$.
62. Show that for any ideal I in a commutative ring R ,

$$\sqrt{I} = \{p \in R : p^k \in I \text{ for some } k \in \mathbb{N}^+\}$$

is a radical ideal.

63. Prove the first five properties of the operator \mathcal{Z} listed in §13.12.
64. Let I and J be ideals of a commutative ring R . Let IJ be the set of finite sums $i_1j_1 + \dots + i_mj_m$ where $m \in \mathbb{N}^+$, each $i_s \in I$, and each $j_s \in J$. (a) Check that IJ is an ideal of R , and $IJ \subseteq I \cap J$. (b) Prove that $\mathcal{Z}(IJ) = \mathcal{Z}(I \cap J) = \mathcal{Z}(I) \cup \mathcal{Z}(J)$.
65. Prove the first three properties of the operator \mathcal{A} listed in §13.12.
66. True or false? Explain each answer. (a) For all $x \in V$, there exists $f \in V^*$ with $f(x) \neq 0$. (b) For all vector spaces Z , $\dim(Z) = \dim(Z^*)$. (c) Given $\dim(V) = n$, every nonzero $f \in V^*$ satisfies $\dim(\ker(f)) = n - 1$. (d) For all vector spaces Z , $\text{ev} : Z \rightarrow Z^{**}$ is an isomorphism. (e) For all $S, T \subseteq V^*$, if $\mathcal{Z}(S) = \mathcal{Z}(T)$ then $S = T$. (f) For all $x \neq y$ in V , there exists $f \in V^*$ with $f(x) \neq f(y)$. (g) For all linearly independent lists (f_1, \dots, f_k) in V^* , there exist $x_1, \dots, x_k \in V$ with $f_i(x_j) = \delta_{i,j}$ for $1 \leq i, j \leq k$. (h) For all subspaces U, Z of a real inner product space V , if $U^\perp = Z^\perp$ then $U = Z$. (i) The map $D : \text{Hom}_F(V, W) \rightarrow \text{Hom}_F(W^*, V^*)$ given by $D(T) = T^*$ is F -linear. (j) For all linear maps $S, T : V \rightarrow V$, $(S \circ T)^* = S^* \circ T^*$. (k) For all $T \subseteq V^*$, $\mathcal{Z}(\mathcal{A}(\mathcal{Z}(T))) = \mathcal{Z}(T)$.

Metric Spaces and Hilbert Spaces

So far, our study of linear algebra has focused mostly on *finite-dimensional* vector spaces and the linear maps between such spaces. In this chapter, we want to give the reader a small taste of the ideas that are needed to treat linear algebra in the infinite-dimensional setting. Infinite-dimensional vector spaces arise frequently in analysis, where one studies vector spaces of functions. Not surprisingly, to understand such spaces in detail, one needs not just the algebraic concepts of linear algebra but also some tools from analysis and topology. In particular, the notions of the length of a vector, the distance between two vectors, and the convergence of a sequence of vectors are key ingredients in the study of infinite-dimensional spaces.

The first part of this chapter develops the analytical concepts we will need in the context of *metric spaces*, which are sets in which one has defined the distance between any two points. The distance function (also called a *metric*) satisfies a few basic axioms, from which many fundamental properties of convergent sequences can be derived. We use convergent sequences to define other topological concepts such as closed sets, open sets, continuous functions, compact sets, and complete spaces. Our discussion of metric spaces is far from comprehensive, but it does provide a self-contained account of the analytic material needed for our treatment of Hilbert spaces in the second part of the chapter.

Intuitively, a Hilbert space is a complex vector space (often infinite-dimensional) in which we can define both the *length* of a vector and the notion of *orthogonality* of vectors, which generalizes the geometric idea of perpendicular vectors. Orthogonality is defined via a scalar product similar to the dot product in \mathbb{R}^3 , but using complex scalars. A crucial technical condition in the definition of a Hilbert space is the requirement that it be *complete* as a metric space (which means, informally, that if the terms in a given sequence get arbitrarily close to each other, then the sequence must actually converge to some point in the space).

We will see that the assumption of completeness provides an analytic substitute for the dimension-counting arguments that one often needs when proving facts about finite-dimensional inner product spaces. For example, given any subspace W of a finite-dimensional inner product space V , the theorem that $V = W \oplus W^\perp$ (see §13.9) used dimension counting in its proof. The corresponding result in Hilbert spaces requires an additional hypothesis (the assumption that the subspace W be a *closed* set), and the proof uses completeness in an essential way.

Orthonormal bases play a central role in the theory of finite-dimensional inner product spaces. The analogous concept for infinite-dimensional Hilbert spaces is the idea of a *maximal orthonormal set*. We will see that every Hilbert space has a maximal orthonormal set X , and every vector in the space can be expressed as an “infinite linear combination” of the vectors in X . (One needs analytical ideas to give a precise definition of what this means.) These results lead to a classification theorem showing that every abstract Hilbert space is isomorphic to a concrete Hilbert space consisting of “square-summable functions” defined on X . After proving these results, the chapter closes with a discussion of spaces of continuous linear maps, the identification of a Hilbert space with its dual space, and adjoint operators.

14.1 Metric Spaces

Metric spaces were discussed briefly in §10.5; we repeat the relevant definitions here to keep this chapter self-contained. A *metric space* is a set X together with a function (called a *metric* or *distance function*) that measures the distance between any two points in the set X . The distance function $d : X \times X \rightarrow \mathbb{R}$ must satisfy the following conditions for all $x, y, z \in X$. First, $0 \leq d(x, y) < \infty$, and $d(x, y) = 0$ iff $x = y$. Second, $d(x, y) = d(y, x)$. Third, $d(x, z) \leq d(x, y) + d(y, z)$. The second axiom is called *symmetry*; the third axiom is called the *triangle inequality*. If several metric spaces are being considered, we sometimes write d_X for the metric defined on X .

What are some examples of metric spaces? The set \mathbb{R} of real numbers is a metric space, with distance function $d(x, y) = |x - y|$ for all $x, y \in \mathbb{R}$. More generally, for all $m \in \mathbb{N}^+$, \mathbb{R}^m is a metric space under the *Euclidean distance function*, defined by setting $d_2(x, y) = \sqrt{\sum_{i=1}^m |x_i - y_i|^2}$ for all $x = (x_1, \dots, x_m)$ and $y = (y_1, \dots, y_m)$ in \mathbb{R}^m . Analogous formulas define metrics on the spaces \mathbb{C} and \mathbb{C}^m . We will verify the metric space axioms for these examples later in the chapter, after we introduce Hilbert spaces.

Here are a few more abstract examples of metric spaces. Given any set X and $x, y \in X$, define $d(x, y) = 0$ if $x = y$, and $d(x, y) = 1$ if $x \neq y$. The first two axioms for a metric are immediately verified. To see that the triangle inequality must hold, suppose that it failed. Then there would exist $x, y, z \in X$ with $d(x, z) > d(x, y) + d(y, z)$. Since the distance function only takes values 0 and 1, the inequality just written can only occur if $d(x, z) = 1$ and $d(x, y) = d(y, z) = 0$. But then the definition of d gives $x = y = z$, hence $x = z$, which contradicts $d(x, z) = 1$. This distance function d is called the *discrete metric* on X .

Given any metric space (X, d) and any subset S of X , the set S becomes a metric space by restricting the domain of the distance function d from $X \times X$ to $S \times S$. We call S a *subspace* of the metric space X .

Next, suppose X and Y are metric spaces with metrics d_X and d_Y . Consider the *product space* $Z = X \times Y$, which is the set of ordered pairs (x, y) with $x \in X$ and $y \in Y$. We define the *product metric* $d = d_Z : Z \times Z \rightarrow \mathbb{R}$ by setting

$$d((x_1, y_1), (x_2, y_2)) = d_X(x_1, x_2) + d_Y(y_1, y_2) \text{ for all } (x_1, y_1), (x_2, y_2) \in Z.$$

Let us verify the metric space axioms for Z and d . Fix $z_1 = (x_1, y_1)$, $z_2 = (x_2, y_2)$, and $z_3 = (x_3, y_3)$ in Z , where $x_i \in X$ and $y_i \in Y$ for $i = 1, 2, 3$. First, since $0 \leq d_X(x_1, x_2) < \infty$ and $0 \leq d_Y(y_1, y_2) < \infty$, we have $0 \leq d(z_1, z_2) = d_X(x_1, x_2) + d_Y(y_1, y_2) < \infty$. Furthermore, since the sum of two nonnegative real numbers is zero iff both summands are zero, we have $d(z_1, z_2) = 0$ iff $d_X(x_1, x_2) = 0$ and $d_Y(y_1, y_2) = 0$ iff $x_1 = x_2$ and $y_1 = y_2$ iff $(x_1, y_1) = (x_2, y_2)$ iff $z_1 = z_2$. Second, $d(z_1, z_2) = d_X(x_1, x_2) + d_Y(y_1, y_2) = d_X(x_2, x_1) + d_Y(y_2, y_1) = d(z_2, z_1)$ by symmetry of d_X and d_Y . Third, using the triangle inequality for d_X and d_Y , we compute

$$\begin{aligned} d(z_1, z_3) &= d_X(x_1, x_3) + d_Y(y_1, y_3) \leq d_X(x_1, x_2) + d_X(x_2, x_3) + d_Y(y_1, y_2) + d_Y(y_2, y_3) \\ &= d_X(x_1, x_2) + d_Y(y_1, y_2) + d_X(x_2, x_3) + d_Y(y_2, y_3) = d_Z(z_1, z_2) + d_Z(z_2, z_3). \end{aligned}$$

Similar calculations show that the distance function $d' : Z \times Z \rightarrow \mathbb{R}$ given by $d'((x_1, y_1), (x_2, y_2)) = \max(d_X(x_1, x_2), d_Y(y_1, y_2))$ satisfies the metric space axioms. The metrics d and d' on the set Z are not equal, but we will see in Exercise 9 that they are equivalent for many purposes.

We can iterate the definition of d (or d') to define metrics on products of finitely many

metric spaces X_1, X_2, \dots, X_m . Specifically, let $X = X_1 \times X_2 \times \dots \times X_m$ and let $x_i, y_i \in X_i$ for $1 \leq i \leq m$. The formulas

$$\begin{aligned} d_1((x_1, \dots, x_m), (y_1, \dots, y_m)) &= d_{X_1}(x_1, y_1) + d_{X_2}(x_2, y_2) + \dots + d_{X_m}(x_m, y_m), \\ d_\infty((x_1, \dots, x_m), (y_1, \dots, y_m)) &= \max(d_{X_1}(x_1, y_1), d_{X_2}(x_2, y_2), \dots, d_{X_m}(x_m, y_m)) \end{aligned}$$

both define metrics on the product space X . In particular, these constructions provide two additional metrics on the spaces \mathbb{R}^m and \mathbb{C}^m that are not equal to the Euclidean metric given earlier. Although all three metrics are equivalent in some respects (Exercise 6), we will see later that the Euclidean metric d_2 has additional geometric structure compared to d_1 and d_∞ .

As a final remark, note that the constructions given here do *not* automatically generalize to products of infinitely many metric spaces. The reason is that the sum or maximum of an infinite sequence of positive real numbers may be $+\infty$, which cannot be a value of the distance function, by the first metric space axiom.

14.2 Convergent Sequences

Formally, a *sequence* of points in a set X is a function $\mathbf{x} : \mathbb{N} \rightarrow X$. We usually present such a sequence as an “infinite list” $\mathbf{x} = (x_0, x_1, x_2, \dots, x_n, \dots) = (x_n : n \in \mathbb{N})$, where $x_n = \mathbf{x}(n) \in X$ is the “ n ’th term” of the sequence. We sometimes index a sequence starting at x_1 instead of x_0 , or starting at x_i for any fixed $i \in \mathbb{Z}$. A *subsequence* of $(x_0, x_1, x_2, \dots, x_n, \dots)$ is a sequence of the form $(x_{k_0}, x_{k_1}, x_{k_2}, \dots, x_{k_n}, \dots)$ where $0 \leq k_0 < k_1 < k_2 < \dots < k_n < \dots$ is a strictly increasing sequence of integers.

Given a metric space (X, d) , a sequence (x_n) of points in X , and a point $y \in X$, we say that this sequence *converges to y* in the metric space iff for all $\epsilon > 0$, there exists $n_0 \in \mathbb{N}$ such that for all $n \geq n_0$, $d(x_n, y) < \epsilon$. When this condition holds, we write $x_n \rightarrow y$ or $\lim_{n \rightarrow \infty} x_n = y$ and call y the *limit* of the sequence (x_n) . Intuitively, the formal definition means that for any given positive distance, no matter how small, the sequence (x_n) will eventually come closer than that distance to the limit y and will remain closer than that distance forever. (The phrasing of this intuitive description arises by viewing the subscript n as a time index, so that x_n is the “location of the sequence at time n .”)

Not every sequence (x_n) in a metric space (X, d) will converge. Those sequences that do converge to some point $y \in X$ are called *convergent* sequences. We note that the limit of a sequence, when it exists, must be unique. For suppose a particular sequence (x_n) converged to both y and z in X . We will show that $y = z$. For any fixed $\epsilon > 0$, we can choose $n_1 \in \mathbb{N}$ such that $n \geq n_1$ implies $d(x_n, y) < \epsilon/2$. Similarly, there exists $n_2 \in \mathbb{N}$ such that $n \geq n_2$ implies $d(x_n, z) < \epsilon/2$. Picking $n = \max(n_1, n_2)$, we conclude that $d(y, z) \leq d(y, x_n) + d(x_n, z) = d(x_n, y) + d(x_n, z) < \epsilon/2 + \epsilon/2 = \epsilon$. Thus the fixed number $d(y, z)$ is less than every positive number ϵ , forcing $d(y, z) = 0$ and $y = z$.

Here are some examples of sequences in metric spaces. In any metric space, the constant sequence (y, y, y, \dots) converges to y . In the metric space \mathbb{R} , the sequences $(1/n : n \in \mathbb{N}^+)$, $(1/n^2 : n \in \mathbb{N}^+)$, and $(1/2^n : n \in \mathbb{N})$ all converge to zero. The sequence $(1, -1, 1, -1, 1, -1, \dots) = ((-1)^n : n \geq 0)$ does not converge. However, the subsequence of odd-indexed terms $(-1, -1, -1, \dots)$ converges to -1 , and the subsequence of even-indexed terms $(1, 1, 1, \dots)$ converges to 1 . We see that a subsequence of a non-convergent sequence may converge, and different subsequences might converge to different limits. On the other hand, the sequence $(n : n \in \mathbb{N})$ in \mathbb{R} is non-convergent and has no convergent subsequence (infinity is not allowed as a limit, since it is not in the set \mathbb{R}).

Suppose $(x_n : n \in \mathbb{N})$ is a convergent sequence in a metric space (X, d) with limit y . We now show that *all subsequences of (x_n) also converge to y* . Fix such a subsequence $(x_{k_0}, x_{k_1}, \dots)$. Given $\epsilon > 0$, there exists $m_0 \in \mathbb{N}$ such that for all $m \geq m_0$, $d(x_m, y) < \epsilon$. Since $k_0 < k_1 < \dots$ is a strictly increasing sequence of integers, we can find $n_0 \in \mathbb{N}$ such that $k_n \geq m_0$ for all $n \geq n_0$. (It suffices to take $n_0 = m_0$, although a smaller n_0 may also work.) Now, for any $n \geq n_0$, $d(x_{k_n}, y) < \epsilon$, confirming that $\lim_{n \rightarrow \infty} x_{k_n} = y$. We will develop further connections between the limit of a sequence (if any) and the limits of its subsequences when we discuss compactness and completeness.

14.3 Closed Sets

Let (X, d) be any metric space. A subset C of X is called a *closed set* (relative to the metric d) iff for every sequence $(x_n : n \in \mathbb{N})$ with all terms x_n belonging to C , if x_n converges to some point $y \in X$, then y also belongs to C . Intuitively, the set is called “closed” because we can never go outside the set by passing to the limit of a convergent sequence of points all of which lie in the set. In §1.4, we saw that many algebraic subsystems (subgroups, ideals, etc.) can be defined in terms of certain “closure conditions,” where performing various operations on elements of the set produces another object in that same set. Here, we can say that C is closed (in the sense just defined) iff C is closed under the operation of taking the limit of a convergent sequence of points in C . Negating the definition, note that C is *not closed* iff there exists a sequence (x_n) with all $x_n \in C$ such that $x_n \rightarrow y$ for some $y \in X$, but $y \notin C$.

Here are some examples of closed sets. First, for any fixed $y \in X$, the one-point set $C = \{y\}$ is closed. To see why, note that the only sequence of points in C is the constant sequence (y, y, y, \dots) . This sequence evidently converges to y , which belongs to C , and so C is indeed closed. Second, the empty set \emptyset is a closed subset of X . For if it were not closed, there would be a convergent sequence $(x_n : n \in \mathbb{N})$ with all $x_n \in \emptyset$ converging to a limit $y \notin \emptyset$. But the condition $x_n \in \emptyset$ is impossible, so there can be no such sequence. Third, the entire space X is a closed subset of X . For, given points $x_n \in X$ such that $x_n \rightarrow y \in X$, we certainly have $y \in X$, so that X is closed.

We continue by giving examples of closed sets and non-closed sets in the metric space \mathbb{R} . For any fixed $a \leq b$, the *closed interval* $[a, b] = \{x \in \mathbb{R} : a \leq x \leq b\}$ is a closed set, as the name suggests. We prove this by contradiction. If $[a, b]$ is not closed, choose a sequence (x_n) with all $x_n \in [a, b]$, such that (x_n) converges to some real number $y \notin [a, b]$. In the case $y < a$, note that $\epsilon = a - y > 0$. For this ϵ , there is $n_0 \in \mathbb{N}$ such that $n \geq n_0$ implies $d(x_n, y) < \epsilon$. In particular, x_{n_0} must satisfy $y - \epsilon < x_{n_0} < y + \epsilon = a$, contradicting $x_{n_0} \in [a, b]$. Similarly, the case $y > b$ leads to a contradiction. So $[a, b]$ is closed. On the other hand, for fixed $a < b$, the open interval $(a, b) = \{x \in \mathbb{R} : a < x < b\}$ is not closed. For, one can check that the sequence defined by $x_n = a + (b - a)/(2n)$ for $n \geq 1$ has all $x_n \in (a, b)$ and converges to $a \notin (a, b)$. Similarly, the half-open intervals $(a, b]$ and $[a, b)$ are not closed. The set \mathbb{Z} of integers is a closed subset of \mathbb{R} . We again argue by contradiction. If \mathbb{Z} is not closed, choose a convergent sequence (x_n) of integers with $x_n \rightarrow y$, where $y \in \mathbb{R}$ is not an integer. We know y lies between two consecutive integers, say $k < y < k + 1$. Take $\epsilon = \min(y - k, k + 1 - y) > 0$. For large enough n , we must have $y - \epsilon < x_n < y + \epsilon$. But no real number in this range is an integer, which contradicts $x_n \in \mathbb{Z}$.

Let us return to the case of a general metric space (X, d) . Given any (possibly infinite) collection $\{C_i : i \in I\}$ of closed sets in X , we claim the intersection $C = \bigcap_{i \in I} C_i$ is closed. For, suppose (x_n) is a sequence in C converging to some $y \in X$; we must prove $y \in C$. Fix an index $i \in I$. Since all x_n lie in C , we know all $x_n \in C_i$. Because C_i is a closed set, it

follows that the limit y belongs to C_i . This holds for every i , so y belongs to the intersection C of all the C_i .

Next we show that if C and D are closed sets in X , then the union $C \cup D$ is also closed. Let (x_n) be a sequence of points in $C \cup D$ converging to a limit $y \in X$; we must prove $y \in C \cup D$. Consider two cases. Case 1: there are only finitely many indices n with $x_n \in C$. Let n_0 be the largest index with $x_{n_0} \in C$, or $n_0 = 0$ if every x_n is in D . The subsequence $(x_{n_0+1}, x_{n_0+2}, \dots)$ of the original sequence still converges to y , and all points in this subsequence belong to D . As D is closed, we must have $y \in D$, so that $y \in C \cup D$ as well. Case 2: there are infinitely many indices $k_0 < k_1 < k_2 < \dots$ with $x_{k_n} \in C$. Then the subsequence $(x_{k_0}, x_{k_1}, \dots)$ of the original sequence still converges to y , and all points in this subsequence belong to C . As C is closed, we must have $y \in C$, so that $y \in C \cup D$.

It follows by induction that if $m \in \mathbb{N}^+$ and C_1, C_2, \dots, C_m are closed sets in (X, d) , then $C_1 \cup C_2 \cup \dots \cup C_m$ is also closed. In particular, since one-point sets are closed, we see that *all finite subsets of X are closed*. However, the union of an infinite collection of closed subsets may or may not be closed. On one hand, we saw that $\mathbb{Z} = \bigcup_{n \in \mathbb{Z}} \{n\}$ is a closed set in \mathbb{R} . On the other hand, the union C of the one-point sets $\{1/n\}$ for $n \in \mathbb{N}^+$ is not closed in \mathbb{R} , since $(1/n : n \in \mathbb{N}^+)$ is a sequence in C converging to the point $0 \notin C$.

To summarize: *in any metric space X , \emptyset and X are closed. Finite subsets of X are closed. The union of finitely many closed sets is closed. The intersection of arbitrarily many closed sets is closed.*

14.4 Open Sets

Given a point x in a metric space (X, d) and a real $r > 0$, the *open ball of radius r and center x* is $B(x; r) = \{y \in X : d(x, y) < r\}$. For example, in \mathbb{R} , $B(x; r)$ is the open interval $(x - r, x + r)$. In \mathbb{C} or \mathbb{R}^2 with the Euclidean metric, $B(x; r)$ is the interior of a circle with center x and radius r . In a discrete metric space X , $B(x; r) = \{x\}$ for all $r \leq 1$, whereas $B(x; r) = X$ for all $r > 1$.

A subset U of a metric space (X, d) is called an *open set* (relative to the metric d) iff for every $x \in U$, there exists $\epsilon > 0$ (depending on x) such that $B(x; \epsilon) \subseteq U$. Intuitively, the set U is called open because all points sufficiently close to a point in U are also in U .

Here are some examples of open sets in general metric spaces. First, *every open ball is an open set*, as the name suggests. To prove this, consider an open ball $U = B(y; r)$ and fix some $x \in U$. We know $d(y, x) < r$, so the number $\epsilon = r - d(y, x)$ is strictly positive. We will show $B(x; \epsilon) \subseteq U$. Fix $z \in B(x; \epsilon)$; is $z \in U$? We know $d(x, z) < \epsilon$, so the triangle inequality gives $d(y, z) \leq d(y, x) + d(x, z) < d(y, x) + \epsilon = r$. This shows that $z \in B(y; r) = U$, as needed.

Second, *the entire space X is an open subset of X* . For, given $x \in X$, we can take $\epsilon = 1$ and note that $B(x; \epsilon)$ is a subset of X by definition.

Third, *the empty set \emptyset is an open subset of X* . For if \emptyset were not open, there must exist $x \in \emptyset$ such that for all $\epsilon > 0$, $B(x; \epsilon)$ is not a subset of \emptyset . But the existence of $x \in \emptyset$ is impossible.

Fourth, *given any collection $\{U_i : i \in I\}$ of open subsets of X , the union $U = \bigcup_{i \in I} U_i$ is also open in X* . To prove this, fix x in the union U of the U_i . We know $x \in U_i$ for some $i \in I$. Since U_i is open, there exists $\epsilon > 0$ with $B(x; \epsilon) \subseteq U_i$. As U_i is a subset of U , we also have $B(x; \epsilon) \subseteq U$, so U is open.

Fifth, *if U and V are open subsets of X , then $U \cap V$ is open in X* . For the proof, fix $x \in U \cap V$. Since $x \in U$, there is $\epsilon_1 > 0$ with $B(x; \epsilon_1) \subseteq U$. Since $x \in V$, there is $\epsilon_2 > 0$

with $B(x; \epsilon_2) \subseteq V$. For $\epsilon = \min(\epsilon_1, \epsilon_2) > 0$, we see that $B(x; \epsilon)$ is contained in both U and V , hence is a subset of $U \cap V$. Thus, $U \cap V$ is open. By induction, it follows that if $m \in \mathbb{N}^+$ and U_1, \dots, U_m are open subsets of X , then $U_1 \cap \dots \cap U_m$ is also open. However, the intersection of infinitely many open subsets of X need not be open. For instance, all the sets $B(0; 1/n) = (-1/n, 1/n)$ are open subsets of \mathbb{R} (being open balls). But their intersection, namely $\{0\}$, is not open in \mathbb{R} , since for every $\epsilon > 0$, $B(0; \epsilon) = (-\epsilon, \epsilon)$ is not a subset of $\{0\}$.

Open sets and closed sets are related in the following way: *a set U is open in X iff the complement $C = X \sim U$ is closed in X* . We prove the contrapositive in both directions. First assume $C = X \sim U$ is not closed in X . Then there is a sequence (x_n) of points of C and a point $x \in X$ such that $x_n \rightarrow x$ but $x \notin C$. Note x is a point of $X \sim C = U$. For every $\epsilon > 0$, there is $n_0 \in \mathbb{N}$ such that for all $n \geq n_0$, $d(x_n, x) < \epsilon$. Each such x_n lies in both $C = X \sim U$ and $B(x; \epsilon)$, which shows that $B(x; \epsilon)$ cannot be a subset of U . As ϵ was arbitrary, we see that U is not open.

Conversely, assume U is not open; we prove $C = X \sim U$ is not closed. There is a point $x \in U$ such that for all $\epsilon > 0$, $B(x; \epsilon)$ is not a subset of U . Taking $\epsilon = 1/n$ for each $n \in \mathbb{N}^+$, we obtain points x_n with $x_n \in B(x; 1/n)$ but $x_n \notin U$. Thus, (x_n) is a sequence of points in C . Furthermore, we claim x_n converges to x in X . To see this, fix $\epsilon > 0$, and choose $n_0 \in \mathbb{N}$ with $1/n_0 < \epsilon$. For all $n \geq n_0$, $d(x_n, x) < 1/n \leq 1/n_0 < \epsilon$, as needed. Since the limit of (x_n) is $x \notin C$, we see that C is not closed.

We remark that most subsets of most metric spaces are neither open nor closed. For instance, half-open intervals $(a, b]$ in \mathbb{R} are neither open nor closed. It is possible for a subset of a metric space to be both open and closed; consider \emptyset and X , for example.

14.5 Continuous Functions

The concept of *continuity* plays a central role in calculus. This concept can be generalized to the setting of metric spaces as follows. Let (X, d_X) and (Y, d_Y) be two metric spaces. A function $f : X \rightarrow Y$ is *continuous on X* iff whenever (x_n) is a sequence of points in X converging to some $x \in X$, the sequence $(f(x_n))$ in Y converges to $f(x)$. Briefly, continuity of f means that whenever $x_n \rightarrow x$ in X , $f(x_n) \rightarrow f(x)$ in Y . Using limit notation, continuity of f means

$$f\left(\lim_{n \rightarrow \infty} x_n\right) = \lim_{n \rightarrow \infty} f(x_n),$$

so that f “commutes” with the operation of taking the limit of a convergent sequence. We say f is *continuous at the point x^* of X* iff whenever $x_n \rightarrow x^*$ in X , $f(x_n) \rightarrow f(x^*)$ in Y .

One readily checks that a constant function ($f(x) = y_0$ for all $x \in X$) is continuous, as is the identity function $\text{id}_X : X \rightarrow X$. We will see later that addition and multiplication (viewed as functions from the product metric space $\mathbb{R} \times \mathbb{R}$ to \mathbb{R}) are continuous. Let us check that $f : \mathbb{R} \rightarrow \mathbb{R}$, given by $f(x) = -3x$ for $x \in \mathbb{R}$, is continuous. Suppose $x_n \rightarrow x$ in \mathbb{R} ; does $-3x_n \rightarrow -3x$? Fix $\epsilon > 0$, and choose n_0 so that $n \geq n_0$ implies $d(x_n, x) < \epsilon/3$. Now notice that $d(-3x_n, -3x) = |(-3x_n) - (-3x)| = 3|x_n - x| = 3d(x_n, x)$. So, $n \geq n_0$ implies $d(f(x_n), f(x)) < \epsilon$, as needed.

A fundamental fact about continuity is that *compositions of continuous functions are continuous*. In detail, let $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ be continuous functions between metric spaces X , Y , and Z ; we show $g \circ f : X \rightarrow Z$ is continuous. To do so, let (x_n) be a sequence in X converging to $x \in X$. By continuity of f , the sequence $(f(x_n))$ converges to $f(x)$ in Y . By continuity of g , the sequence $(g(f(x_n)))$ converges to $g(f(x))$ in Z . So, the sequence $((g \circ f)(x_n))$ converges to $(g \circ f)(x)$ in Z , proving continuity of $g \circ f$.

A more subtle, but equally fundamental, property of continuity is that $f : X \rightarrow Y$ is continuous iff for every closed set D in Y , the inverse image $f^{-1}[D] = \{x \in X : f(x) \in D\}$ is closed in X . To prove the forward direction, assume f is continuous, fix a closed set D in Y , and consider a sequence (x_n) of points in $f^{-1}[D]$ that converge to some limit $x^* \in X$. To see that $f^{-1}[D]$ is closed, we must prove $x^* \in f^{-1}[D]$. Now, each x_n is in $f^{-1}[D]$, so $f(x_n) \in D$ for all $n \in \mathbb{N}$. As x_n converges to x^* , continuity of f tells us that $f(x_n)$ converges to $f(x^*)$. As D was assumed to be closed in Y , we deduce $f(x^*) \in D$, and hence $x^* \in f^{-1}[D]$ as needed.

For the other direction, assume that $f : X \rightarrow Y$ is not continuous. So there is a sequence (x_n) in X converging to a point $x \in X$, such that $f(x_n)$ does not converge to $y = f(x)$ in Y . This means that there exists $\epsilon > 0$ such that for every $n_0 \in \mathbb{N}$ there exists $n \geq n_0$ with $d(f(x_n), y) \geq \epsilon$. It follows that we can find a subsequence (x'_n) of (x_n) , consisting of all terms x_n such that $d(f(x_n), f(x)) \geq \epsilon$. We know (x'_n) still converges to x . Note that $D = \{z \in Y : d_Y(z, y) \geq \epsilon\}$ is a closed set in Y , being the complement of the open ball $B(y; \epsilon)$. To complete the proof, we show $f^{-1}[D]$ is not a closed set in X . Note that each term x'_n lies in $f^{-1}[D]$, since $f(x'_n) \in D$ by construction. But the limit x of the sequence (x'_n) does not lie in $f^{-1}[D]$, because $f(x) = y$ satisfies $d_Y(y, y) = 0$, so that $f(x) \notin D$.

Since open sets are complements of closed sets, we readily deduce the following characterization of continuous functions: $f : X \rightarrow Y$ is continuous iff for every open set V in Y , the inverse image $f^{-1}[V]$ is open in X . For suppose f is continuous and V is open in Y . Then $Y \sim V$ is closed in Y , so $f^{-1}[Y \sim V]$ is closed in X . From set theory, we know $f^{-1}[Y \sim V] = f^{-1}[Y] \sim f^{-1}[V] = X \sim f^{-1}[V]$. Therefore, $f^{-1}[V]$ is the complement of a closed set in X , so it is open in X . The converse is proved similarly. This characterization of continuity is closely related to the definition often given in calculus using ϵ 's and δ 's. Specifically, $f : X \rightarrow Y$ is continuous iff for all $x_0 \in X$ and $\epsilon > 0$, there exists $\delta > 0$ such that for all $x \in X$, if $d_X(x, x_0) < \delta$, then $d_Y(f(x), f(x_0)) < \epsilon$. We ask the reader to prove the equivalence of this condition to the condition involving open sets in Exercise 36.

14.6 Compact Sets

We know that not every sequence in a general metric space is convergent. It would be convenient if we could take a non-convergent sequence (x_n) and extract a convergent subsequence from it. However, even this is not always possible. For example, in \mathbb{R} , no subsequence of the sequence $(n : n \in \mathbb{N})$ converges. On the other hand, if we insist that all terms of the sequence (x_n) come from a closed interval $[a, b]$, it can be shown that (x_n) must have a convergent subsequence (Exercise 41). This suggests that we may have to restrict our attention to sequences coming from a sufficiently nice subset of the metric space. In a general metric space (X, d) , a subset K of X is called (*sequentially*) compact iff for every sequence (x_n) of points in K , there exists a subsequence (x_{k_n}) and a point y belonging to K such that $x_{k_n} \rightarrow y$.

Every finite subset of X is compact. For suppose $K = \{y_1, \dots, y_m\}$ is nonempty and finite, and (x_n) is any sequence of points in K . At least one y_j must occur infinitely often in the sequence (x_n) . So this sequence has a constant subsequence (y_j, y_j, y_j, \dots) , which evidently converges to the point $y_j \in K$. One checks by negating the definition that \emptyset is also compact.

Every compact subset of X must be closed in X . For, suppose $K \subseteq X$ is compact, (x_n) is a sequence with all $x_n \in K$, and x_n converges to a point x in X . On one hand, every subsequence of (x_n) converges to the unique limit x . On the other hand, the definition of

compactness shows that some subsequence of (x_n) converges to a point of K . Therefore, x must belong to K , proving that K is closed. One may show that *unions of finitely many compact sets are compact, whereas intersections of arbitrarily many compact sets are compact*. The proof (Exercise 38) imitates the proofs of the analogous properties of closed sets. Similarly, *a closed subset of a compact set is compact*.

We say that a subset S of a nonempty metric space (X, d) is *bounded* iff there exists $z \in X$ and $M \in \mathbb{R}$ such that for all $x \in S$, $d(z, x) \leq M$. *Every compact subset of X must be bounded.* For, suppose $S \subseteq X$ is not bounded. We will construct a sequence (x_n) of points in S that has no convergent subsequence. Fix $z \in X$. For each $n \in \mathbb{N}$, we can find $x_n \in S$ such that $d(z, x_n) > n$. To get a contradiction, suppose some subsequence (x_{k_n}) converged to some point $y \in X$. Choose $n_0 \in \mathbb{N}$ so that $n \geq n_0$ implies $d(y, x_{k_n}) < 1$. Then choose an integer $n \geq n_0$ with $k_n \geq d(z, y) + 1$. For this n , $d(z, x_{k_n}) \leq d(z, y) + d(y, x_{k_n}) < d(z, y) + 1 \leq k_n$, which contradicts the choice of x_{k_n} .

So far, we have seen that *every compact subset of a metric space is closed and bounded*. In \mathbb{R}^m and \mathbb{C}^m with any of the metrics d_1 , d_2 , and d_∞ discussed earlier, the converse also holds: *a subset K of \mathbb{R}^m or \mathbb{C}^m is compact iff K is closed and bounded*. The proof of this, which is rather difficult, is sketched in the exercises and can be found in texts on advanced calculus. On the other hand, *in general metric spaces, there can exist closed and bounded sets that are not compact*. For instance, consider $X = \mathbb{Z}$ with the discrete metric. The entire space X is bounded, since $d(0, x) \leq 1$ for all $x \in \mathbb{Z}$. X is also closed in X . But X is not compact, since one sees readily that $(n : n \in \mathbb{N})$ is a sequence in X with no convergent subsequence.

Continuous functions “preserve” compact sets, in the following sense. *If $f : X \rightarrow Y$ is a continuous function between two metric spaces and $K \subseteq X$ is compact in X , then the direct image $f[K] = \{f(x) : x \in K\}$ is compact in Y .* To prove this, assume $f : X \rightarrow Y$ is continuous and K is a compact subset of X . To prove $f[K]$ is compact, let (y_n) be any sequence of points in $f[K]$. Each y_n has the form $y_n = f(x_n)$ for some $x_n \in K$. Now (x_n) is a sequence of points in the compact set K , so there is a subsequence (x_{k_n}) converging to some point $x^* \in K$. By continuity, $(y_{k_n}) = (f(x_{k_n}))$ is a subsequence of (y_n) converging to $f(x^*) \in f[K]$. So $f[K]$ is indeed compact.

Compactness is often defined in terms of open sets. Given any set S in a metric space (X, d) , an *open cover* of S is a collection $\{U_i : i \in I\}$ of open subsets of X such that $S \subseteq \bigcup_{i \in I} U_i$. A *finite subcover* is a finite subcollection $\{U_{i_1}, \dots, U_{i_m}\}$ of the given open cover such that $S \subseteq \bigcup_{j=1}^m U_{i_j}$. Such a finite subcover need not exist. We say that a subset K of X is (*topologically*) *compact* iff for each open cover of K , there does exist a finite subcover. It can be shown that *in metric spaces, sequential compactness is equivalent to topological compactness*. The topological definition applies in more general situations, but our study of Hilbert spaces will only need the more intuitive definition in terms of subsequences.

14.7 Completeness

A *Cauchy sequence* in a metric space (X, d) is a sequence $(x_n : n \in \mathbb{N})$ with the following property: for all $\epsilon > 0$, there exists $n_0 \in \mathbb{N}$ such that for all $m, n \geq n_0$, $d(x_n, x_m) < \epsilon$. This definition says that the terms of a Cauchy sequence get arbitrarily close to each other if we go far enough out in the sequence. In contrast, for a convergent sequence with limit x , the terms of the sequence get arbitrarily close to the limiting point x . Let us explore the relationship between these concepts.

On one hand, *every convergent sequence in a metric space is a Cauchy sequence*. For

suppose (x_n) is a sequence in (X, d) converging to $x \in X$. To prove that this sequence is a Cauchy sequence, fix $\epsilon > 0$. Choose $n_0 \in \mathbb{N}$ so that for all $n \geq n_0$, $d(x_n, x) < \epsilon/2$. Then for all $n, m \geq n_0$, the triangle inequality gives $d(x_n, x_m) \leq d(x_n, x) + d(x, x_m) < \epsilon/2 + \epsilon/2 = \epsilon$.

On the other hand, *a Cauchy sequence in a general metric space may not converge*. For example, consider the set \mathbb{R}^+ of strictly positive real numbers with the metric $d(x, y) = |x - y|$ for $x, y \in \mathbb{R}^+$. The sequence $(1/n : n \in \mathbb{N}^+)$ converges in the larger metric space \mathbb{R} to the unique limit zero, so this sequence is a Cauchy sequence (in \mathbb{R} and in \mathbb{R}^+). But the sequence does not converge to any point of \mathbb{R}^+ . Similarly, consider the metric space \mathbb{Q} of rational numbers with $d(x, y) = |x - y|$ for $x, y \in \mathbb{Q}$. We can find a sequence of rational numbers $(x_n : n \in \mathbb{N})$ converging to the irrational real number $\sqrt{2}$ (for instance, let x_n consist of the decimal expansion of $\sqrt{2}$ truncated n places after the decimal). This sequence is a Cauchy sequence in \mathbb{Q} that does not converge in \mathbb{Q} .

A metric space (X, d) is called *complete* iff every Cauchy sequence (x_n) in X does converge to some point in the space X . The above remarks show that \mathbb{R}^+ and \mathbb{Q} are not complete. On the other hand, *any compact metric space is complete*. For suppose (x_n) is a Cauchy sequence in a compact metric space (X, d) . By compactness, this sequence has a subsequence (x_{k_n}) converging to some $x \in X$. We now argue that the full sequence must also converge to x . Given $\epsilon > 0$, choose $n_0 \in \mathbb{N}$ so that for all $n \geq n_0$, $d(x_{k_n}, x) < \epsilon/2$. Also choose $n_1 \in \mathbb{N}$ so that for all $i, j \geq n_1$, $d(x_i, x_j) < \epsilon/2$. Fix any $i \geq n_1$. We can choose an integer $j \geq n_1$ such that $j = k_n$ for some $n \geq n_0$. Using this j , we see that $d(x_i, x) \leq d(x_i, x_j) + d(x_j, x) < \epsilon/2 + d(x_{k_n}, x) < \epsilon$. This proves that (x_n) does converge to x .

A subset of a complete metric space is complete iff it is closed. For suppose (X, d) is complete, $C \subseteq X$ is closed, and (x_n) is a Cauchy sequence with all $x_n \in C$. Then (x_n) is a Cauchy sequence in X , hence converges to some $x \in X$. Since C is closed, the limit x must actually lie in C . So C is complete. Conversely, if $C \subseteq X$ is not closed, there exists a sequence (x_n) with all $x_n \in C$ converging to some $x \in X \setminus C$. The convergent sequence (x_n) is a Cauchy sequence, and x is its unique limit. Since this limit does not lie in C , (x_n) is a Cauchy sequence in C that does not converge in C . So C is not complete.

Each closed interval $[a, b]$ is a compact, hence complete, subset of \mathbb{R} (see Exercise 41). We use this fact to show that \mathbb{R} is complete. The proof requires the lemma that *a Cauchy sequence in any metric space is bounded* (Exercise 13). Given a Cauchy sequence (x_n) in \mathbb{R} , we can therefore choose $M \in \mathbb{R}$ such that every term x_n lies in the closed interval $[-M, M]$. The completeness of this interval ensures that the given Cauchy sequence converges to some real number.

More generally, \mathbb{R}^k (and similarly \mathbb{C}^k) with the Euclidean metric is complete. To prove this, let (x_n) be a Cauchy sequence in \mathbb{R}^k , where $x_n = (x_n(1), x_n(2), \dots, x_n(k))$ for certain real numbers $x_n(i)$. For $1 \leq i \leq k$, the inequality $|x_n(i) - x_m(i)| \leq d_2(x_n, x_m)$ shows that $(x_n(i) : n \in \mathbb{N})$ is a Cauchy sequence in \mathbb{R} . By completeness of \mathbb{R} , this sequence converges to some real number $x(i)$. Let $x = (x(1), x(2), \dots, x(k)) \in \mathbb{R}^k$. Given $\epsilon > 0$, choose integers n_1, \dots, n_k such that $n \geq n_i$ implies $|x_n(i) - x_i| < \epsilon/\sqrt{k}$. Then for $n \geq \max(n_1, \dots, n_k)$, we have

$$d_2(x_n, x) = \sqrt{\sum_{i=1}^k |x_n(i) - x(i)|^2} \leq \sqrt{\sum_{i=1}^k \epsilon^2/k} = \epsilon.$$

So the given Cauchy sequence (x_n) in \mathbb{R}^k converges to x .

14.8 Definition of a Hilbert Space

Having covered the necessary background on metric spaces, we are now ready to define Hilbert spaces. Briefly, a *Hilbert space* is a complex inner product space that is complete relative to the metric induced by the inner product. Let us spell this definition out in more detail.

We begin with the algebraic ingredients of a Hilbert space. The Hilbert space consists of a set H of vectors, together with operations of vector addition $+ : H \times H \rightarrow H$ and scalar multiplication $\cdot : \mathbb{C} \times H \rightarrow H$ satisfying the vector space axioms listed in Table 1.4. There is also defined on H a complex inner product $B : H \times H \rightarrow H$, denoted $B(v, w) = \langle v, w \rangle$, satisfying these axioms for all $v, w, z \in H$ and $c \in \mathbb{C}$: $\langle v + w, z \rangle = \langle v, z \rangle + \langle w, z \rangle$; $\langle cv, z \rangle = c\langle v, z \rangle$; $\langle w, v \rangle = \overline{\langle v, w \rangle}$, where the bar denotes complex conjugation; and for $v \neq 0$, $\langle v, v \rangle$ is a strictly positive real number. It follows (as in §13.10) that the inner product is linear in the first argument and conjugate-linear in the second argument. In other words, for all $c_1, \dots, c_m \in \mathbb{C}$ and $v_1, \dots, v_m, w \in H$, we have

$$\langle c_1 v_1 + \cdots + c_m v_m, w \rangle = c_1 \langle v_1, w \rangle + \cdots + c_m \langle v_m, w \rangle;$$

$$\langle w, c_1 v_1 + \cdots + c_m v_m \rangle = \overline{c_1} \langle w, v_1 \rangle + \cdots + \overline{c_m} \langle w, v_m \rangle.$$

Note that $\langle 0, w \rangle = 0 = \langle w, 0 \rangle$ for all $w \in H$. We say that $v, w \in H$ are *orthogonal* iff $\langle v, w \rangle = 0$, which holds iff $\langle w, v \rangle = 0$. The inner product in a Hilbert space generalizes the dot product from \mathbb{R}^2 and \mathbb{R}^3 , and orthogonality generalizes the geometric concept of perpendicularity in \mathbb{R}^2 and \mathbb{R}^3 .

We use the inner product to define the analytic ingredients of the Hilbert space H , namely the length (or norm) of a vector and the distance between two vectors. For all $v \in H$, define the *norm* of v by setting $\|v\| = \sqrt{\langle v, v \rangle}$. For all $v, w \in H$ and $c \in \mathbb{C}$, the following properties hold: $\|v\| \geq 0$; $\|v\| = 0$ iff $v = 0$; $\|cv\| = |c| \cdot \|v\|$; $|\langle v, w \rangle| \leq \|v\| \cdot \|w\|$ (the *Cauchy–Schwarz inequality*); $\|v + w\| \leq \|v\| + \|w\|$ (the *triangle inequality* for norms). The first two properties follow because $\langle v, v \rangle$ is either a positive real number (when $v \neq 0$) or zero (when $v = 0$). The next property is true because

$$\|cv\| = \sqrt{\langle cv, cv \rangle} = \sqrt{c\bar{c}\langle v, v \rangle} = \sqrt{|c|^2 \langle v, v \rangle} = |c| \sqrt{\langle v, v \rangle} = |c| \cdot \|v\|.$$

The Cauchy–Schwarz inequality has a more subtle proof. The inequality holds if $\langle v, w \rangle = 0$, so we may assume $\langle v, w \rangle \neq 0$, hence $v \neq 0$ and $w \neq 0$. We first prove the inequality in the case where $\|v\| = \|w\| = 1$ and $\langle v, w \rangle$ is a positive real number. Then $\langle v, v \rangle = 1 = \langle w, w \rangle$ and $|\langle v, w \rangle| = \langle v, w \rangle = \langle w, v \rangle$. Using the axioms for the complex inner product, we compute

$$0 \leq \langle v - w, v - w \rangle = \langle v, v \rangle - \langle v, w \rangle - \langle w, v \rangle + \langle w, w \rangle. \quad (14.1)$$

Using our current assumptions on v and w , this inequality becomes $0 \leq 1 - 2|\langle v, w \rangle| + 1$, which rearranges to $|\langle v, w \rangle| \leq 1 = \|v\| \cdot \|w\|$. Still assuming $\|v\| = \|w\| = 1$, we next prove the case where $\langle v, w \rangle$ is an arbitrary nonzero complex scalar. In polar form, $\langle v, w \rangle = re^{i\theta}$ for some real $r > 0$ and some $\theta \in [0, 2\pi)$. Let $v_0 = e^{-i\theta}v$; then $\|v_0\| = |e^{-i\theta}| \cdot \|v\| = \|v\| = 1$, $|\langle v_0, w \rangle| = |e^{-i\theta}\langle v, w \rangle| = |\langle v, w \rangle|$, and $\langle v_0, w \rangle = e^{-i\theta}\langle v, w \rangle = r$ is a positive real number. By the case already proved, $|\langle v_0, w \rangle| \leq \|v_0\| \cdot \|w\|$, and therefore $|\langle v, w \rangle| \leq \|v\| \cdot \|w\|$. Finally, we drop the assumption that $\|v\| = \|w\| = 1$. Write $v = cv_1$ and $w = dw_1$, where $c = \|v\|$, $v_1 = c^{-1}v$, $d = \|w\|$, and $w_1 = d^{-1}w$. We see that $\|v_1\| = c^{-1}\|v\| = 1$ and $\|w_1\| = d^{-1}\|w\| = 1$. On one hand, $0 \neq \langle v, w \rangle = \langle cv_1, dw_1 \rangle = |cd|\langle v_1, w_1 \rangle| = cd|\langle v_1, w_1 \rangle|$. On the other hand, the cases already proved now give $|\langle v_1, w_1 \rangle| \leq \|v_1\| \cdot \|w_1\| = 1$. Therefore,

$$|\langle v, w \rangle| = cd|\langle v_1, w_1 \rangle| \leq cd = \|v\| \cdot \|w\|,$$

completing the proof of the Cauchy–Schwarz inequality.

We can now deduce the triangle inequality for norms from the Cauchy–Schwarz inequality. Given $v, w \in H$, compute

$$\|v + w\|^2 = \langle v + w, v + w \rangle = \langle v, v \rangle + \langle v, w \rangle + \langle w, v \rangle + \langle w, w \rangle = \|v\|^2 + \langle v, w \rangle + \overline{\langle v, w \rangle} + \|w\|^2.$$

The two middle terms add up to twice the real part of $\langle v, w \rangle$, which is at most $2|\langle v, w \rangle| \leq 2\|v\| \cdot \|w\|$. Hence,

$$\|v + w\|^2 \leq \|v\|^2 + 2\|v\| \cdot \|w\| + \|w\|^2 = (\|v\| + \|w\|)^2.$$

Taking the positive square root of both sides gives $\|v + w\| \leq \|v\| + \|w\|$, as needed.

We use the norm to define the metric space structure of H . For $v, w \in H$, define the distance between v and w by $d(v, w) = \|v - w\|$. The properties of the norm derived above quickly imply the required axioms for a metric space (as we saw in §10.5). In particular, the triangle inequality for the metric follows from the triangle inequality for the norm, because

$$d(v, z) = \|v - z\| = \|(v - w) + (w - z)\| \leq \|v - w\| + \|w - z\| = d(v, w) + d(w, z)$$

for all $v, w, z \in H$. Like any metric induced from a norm, the metric on a Hilbert space “respects” translations and dilations, in the sense that $d(v + w, z + w) = d(v, z)$ and $d(cx, cw) = |c|d(v, w)$ for all $v, w, z \in H$ and $c \in \mathbb{C}$.

The final topological ingredient in the definition of a Hilbert space is the assumption that H is *complete* relative to the metric given above. So, every Cauchy sequence in H converges to a point of H . Writing out what this means, we see that whenever $(x_n : n \in \mathbb{N})$ is a sequence of vectors in H such that $\lim_{m,n \rightarrow \infty} \|x_m - x_n\| = 0$, there exists a unique $x \in H$ such that $\lim_{n \rightarrow \infty} \|x_n - x\| = 0$.

We remark that we could define a *real* Hilbert space by the same discussion given above, restricting all scalars to come from \mathbb{R} and using a real-valued inner product. In this case, $\langle v, w \rangle = \langle w, v \rangle$, and the inner product is \mathbb{R} -linear (as opposed to conjugate-linear) in both the first and second positions. Henceforth in this chapter, we continue to consider only complex Hilbert spaces, letting the reader make the required modifications to obtain analogous results for real Hilbert spaces.

14.9 Examples of Hilbert Spaces

A basic example of a Hilbert space is the space \mathbb{C}^n of n -tuples $v = (v_1, \dots, v_n)$, where all $v_k \in \mathbb{C}$. The inner product is defined by $\langle v, w \rangle = \sum_{k=1}^n v_k \overline{w_k}$, which can be written $\langle v, w \rangle = w^* v$ if we think of v and w as column vectors. The norm of v is $\|v\| = \sqrt{\sum_{k=1}^n |v_k|^2}$, and the distance between v and w is the Euclidean distance $d_2(v, w)$ discussed in §14.1. The completeness of \mathbb{C}^n under this metric was proved in §14.7. The other axioms of a Hilbert space (namely, that \mathbb{C}^n is a complex vector space and inner product space) may be routinely verified. In \mathbb{C}^n , the Cauchy–Schwarz inequality $|\langle v, w \rangle| \leq \|v\| \cdot \|w\|$ and the triangle inequality $\|v + w\| \leq \|v\| + \|w\|$ for norms translate to the following facts about sums of complex numbers:

$$\left| \sum_{k=1}^n v_k \overline{w_k} \right| \leq \sqrt{\sum_{k=1}^n |v_k|^2} \cdot \sqrt{\sum_{k=1}^n |w_k|^2}; \quad (14.2)$$

$$\sqrt{\sum_{k=1}^n |v_k + w_k|^2} \leq \sqrt{\sum_{k=1}^n |v_k|^2} + \sqrt{\sum_{k=1}^n |w_k|^2}. \quad (14.3)$$

An example of an infinite-dimensional Hilbert space is the space ℓ_2 of all infinite sequences $v = (v_1, v_2, \dots, v_k, \dots) = (v_k : k \geq 1)$ such that $v_k \in \mathbb{C}$ and $\sum_{k=1}^{\infty} |v_k|^2 < \infty$. The inner product is defined by $\langle v, w \rangle = \sum_{k=1}^{\infty} v_k \overline{w_k}$. We will check the Hilbert space axioms for this example in §14.10.

Both of the previous examples are special cases of a general construction to be described shortly. First, we need a technical digression on general infinite summations. Let X be any set, which could be finite, countably infinite, or uncountable. Let $\mathcal{F}(X)$ be the set of all finite subsets of X . Given a nonnegative real number p_k for each $k \in X$, the sum $\sum_{k \in X} p_k$ is defined to be the least upper bound (possibly ∞) of all sums $\sum_{k \in X'} p_k$, as X' ranges over all finite subsets of X . In symbols, $\sum_{k \in X} p_k = \sup_{X' \in \mathcal{F}(X)} \{\sum_{k \in X'} p_k\}$. Given real numbers r_k for $k \in X$, let $p_k = r_k$ if $r_k \geq 0$ and $p_k = 0$ otherwise; and let $n_k = |r_k|$ if $r_k < 0$ and $n_k = 0$ otherwise. Define $\sum_{k \in X} r_k$ to be $\sum_{k \in X} p_k - \sum_{k \in X} n_k$ provided at most one of the latter sums is ∞ . Finally, given complex numbers z_k for $k \in X$, write $z_k = x_k + iy_k$ for $x_k, y_k \in \mathbb{R}$. Define $\sum_{k \in X} z_k = \sum_{k \in X} x_k + i \sum_{k \in X} y_k$ provided both sums on the right side are finite. Properties of summations over the set X are analogous to properties of infinite series. Some of these properties are covered in the exercises.

Now we can describe the general construction for producing Hilbert spaces. Fix a set X , and let $\ell_2(X)$ be the set of all functions $f : X \rightarrow \mathbb{C}$ such that $\sum_{x \in X} |f(x)|^2 < \infty$. We sometimes think of such a function as a generalized sequence or “ X -tuple” $(f(x) : x \in X) = (f_x : x \in X)$. The set $\ell_2(X)$ becomes a Hilbert space under the following operations. Given $f, g \in \ell_2(X)$ and $c \in \mathbb{C}$, define $f + g$ and cf by setting $(f+g)(x) = f(x)+g(x)$ and $(cf)(x) = c(f(x))$ for all $x \in X$. Define $\langle f, g \rangle = \sum_{x \in X} f(x)\overline{g(x)}$. The norm and metric in this Hilbert space are given by $\|f\| = \sqrt{\sum_{x \in X} |f(x)|^2} < \infty$ and $d(f, g) = \sqrt{\sum_{x \in X} |f(x) - g(x)|^2}$. The verification of the Hilbert space axioms, which is somewhat long and technical, is given in §14.10 below. Note that both previous examples are special cases of this construction, since \mathbb{C}^n is $\ell_2(\{1, 2, \dots, n\})$, and ℓ_2 is $\ell_2(\mathbb{N}^+)$.

In turn, the spaces $\ell_2(X)$ are themselves special cases of Hilbert spaces that arise in the theory of Lebesgue integration. It is beyond the scope of this book to discuss this topic in detail, but we allude to a few facts for readers familiar with this theory. Given any measure space X with measure μ , we let $L_2(X, \mu)$ be the set of measurable functions $f : X \rightarrow \mathbb{C}$ such that $\int_X |f|^2 d\mu < \infty$. This is a complex vector space under pointwise operations on functions, and the inner product is defined by $\langle f, g \rangle = \int_X f\overline{g} d\mu$ for $f, g \in L_2(X, \mu)$. It can be proved that $L_2(X, \mu)$ is a Hilbert space with these operations; the fact that this space is complete is a difficult theorem. By taking μ to be counting measure on an arbitrary set X , we obtain the examples $\ell_2(X)$ discussed above.

Finally, given any Hilbert space H , we can form new examples of Hilbert spaces by considering subspaces of H . An arbitrary vector subspace W of H (a subset closed under zero, addition, and scalar multiplication) will automatically satisfy all the Hilbert space axioms except possibly completeness. Since H is complete, we know from §14.7 that the subspace W will be complete iff W is closed. It follows that *closed subspaces of a Hilbert space are also Hilbert spaces, but non-closed subspaces are not complete*. This example illustrates a general theme that, to obtain satisfactory results in infinite-dimensional settings, one often needs to impose a topological condition (in this case, being closed relative to the metric) in addition to algebraic conditions. We remark that *every finite-dimensional subspace of H is automatically closed*; see Exercise 57. But one can give examples of infinite-dimensional subspaces that are not closed (Exercises 52 and 53).

14.10 Proof of the Hilbert Space Axioms for $\ell_2(X)$

Let X be any set. This section gives the proof that $\ell_2(X)$, with the operations defined above, does satisfy all the axioms for a Hilbert space. We divide the proof into three parts: checking the vector space axioms, checking the inner product axioms, and verifying completeness of the metric.

Vector Space Axioms. Recall that ${}^X\mathbb{C}$ denotes the set of *all* functions $f : X \rightarrow \mathbb{C}$; $\ell_2(X)$ is the subset of ${}^X\mathbb{C}$ consisting of those f satisfying $\sum_{x \in X} |f(x)|^2 < \infty$. We saw in §4.2 that ${}^X\mathbb{C}$ is a complex vector space under pointwise operations on functions. Hence, to see that $\ell_2(X)$ is a vector space under the same operations, it suffices to check that $\ell_2(X)$ is a subspace of ${}^X\mathbb{C}$. First, the zero vector in ${}^X\mathbb{C}$, which is the function sending every $x \in X$ to 0, lies in $\ell_2(X)$ because $\sum_{x \in X} 0^2 = 0 < \infty$. Second, for $f \in \ell_2(X)$ and $c \in \mathbb{C}$, it follows readily from the definition of summation over X (Exercise 58) that $\sum_{x \in X} |(cf)(x)|^2 = \sum_{x \in X} |c(f(x))|^2 = |c|^2 \sum_{x \in X} |f(x)|^2 < \infty$, so that $cf \in \ell_2(X)$. Third, fix $f, g \in \ell_2(X)$; we must show $f + g \in \ell_2(X)$. Let $X' = \{x_1, \dots, x_n\} \in \mathcal{F}(X)$ be any finite subset of X . By definition of $\ell_2(X)$, we know that $\|f\|^2 = \sum_{x \in X} |f(x)|^2$ and $\|g\|^2 = \sum_{x \in X} |g(x)|^2$ are finite. Using the known version (14.3) of the triangle inequality for finite lists of complex numbers, we see that

$$\sum_{k=1}^n |f(x_k) + g(x_k)|^2 \leq \left(\sqrt{\sum_{k=1}^n |f(x_k)|^2} + \sqrt{\sum_{k=1}^n |g(x_k)|^2} \right)^2.$$

Now, since the square root function is increasing, the definition of summation over X shows that $\sqrt{\sum_{k=1}^n |f(x_k)|^2} \leq \sqrt{\sum_{x \in X} |f(x)|^2} = \|f\|$, and similarly $\sqrt{\sum_{k=1}^n |g(x_k)|^2} \leq \|g\|$. So

$$\sum_{x \in X'} |f(x) + g(x)|^2 = \sum_{k=1}^n |f(x_k) + g(x_k)|^2 \leq (\|f\| + \|g\|)^2.$$

This calculation shows that the finite number $(\|f\| + \|g\|)^2$, which does not depend on X' , is an upper bound in \mathbb{R} for all the finite sums $\sum_{x \in X'} |f(x) + g(x)|^2$ as X' ranges over $\mathcal{F}(X)$. Thus the least upper bound of the set of these sums is finite, proving that $f + g \in \ell_2(X)$. In fact, our calculation shows that $\|f + g\|^2 \leq (\|f\| + \|g\|)^2$, giving a direct proof of the triangle inequality for norms in $\ell_2(X)$.

Inner Product Axioms. The main technical point to be checked when verifying the inner product axioms for $\ell_2(X)$ is that $\langle f, g \rangle = \sum_{x \in X} f(x)\overline{g(x)}$ is a well-defined complex number for all $f, g \in \ell_2(X)$. We first prove this under the additional assumption that $f(x)$ and $g(x)$ are *nonnegative real numbers* for all $x \in X$. Given any $X' = \{x_1, \dots, x_n\} \in \mathcal{F}(X)$, the Cauchy–Schwarz inequality for n -tuples (see (14.2)) tells us that

$$\sum_{k=1}^n f(x_k)g(x_k) \leq \sqrt{\sum_{k=1}^n f(x_k)^2} \cdot \sqrt{\sum_{k=1}^n g(x_k)^2}.$$

The right side is at most $\|f\| \cdot \|g\|$, which is a finite upper bound on $\sum_{x \in X'} f(x)g(x)$ that is independent of X' . So $\sum_{x \in X} f(x)g(x) < \infty$ in this case; in fact, the sum is bounded by $\|f\| \cdot \|g\|$.

We next consider the case where $f, g \in \ell_2(X)$ are *real-valued*. Note that $|f|$ has the same squared norm as f , namely $\sum_{x \in X} |f(x)|^2 < \infty$, so $|f|$ (and similarly $|g|$) lie in $\ell_2(X)$. By the case already considered, $\sum_{x \in X} |f(x)g(x)| < \infty$. Now, using a general

property of sums over X (Exercise 58), we can conclude that $\sum_{x \in X} f(x)g(x)$ is a finite real number, and in fact $|\sum_{x \in X} f(x)g(x)| \leq \sum_{x \in X} |f(x)g(x)| \leq \|f\| \cdot \|g\| < \infty$. Finally, for complex-valued $f, g \in \ell_2(X)$, we can write $f = t + iu$ and $g = v + iw$ where $t, u, v, w : X \rightarrow \mathbb{R}$ give the real and imaginary parts of f and g . Since $|t(x)| \leq |f(x)|$ for all $x \in X$, one sees that t (and similarly u, v, w) are real-valued functions in $\ell_2(X)$. Also, $f(x)\bar{g}(x) = [t(x)v(x) + u(x)w(x)] + i[u(x)v(x) - t(x)w(x)]$ for all $x \in X$. By cases already considered, $\sum_{x \in X} t(x)v(x)$, $\sum_{x \in X} u(x)w(x)$, $\sum_{x \in X} u(x)v(x)$, and $\sum_{x \in X} t(x)w(x)$ are all finite. It follows that $\sum_{x \in X} f(x)\bar{g}(x)$ is well-defined. In fact, using Exercise 58, we have the bound $|\sum_{x \in X} f(x)\bar{g}(x)| \leq \sum_{x \in X} |f(x)g(x)| \leq \|f\| \cdot \|g\|$. This gives a direct proof of the Cauchy-Schwarz inequality for $\ell_2(X)$.

Knowing that $\langle f, g \rangle$ is a well-defined complex number, we can prove the axioms for the inner product without difficulty. We prove $\langle f + g, h \rangle = \langle f, h \rangle + \langle g, h \rangle$ for all $f, g, h \in \ell_2(X)$, leaving the remaining axioms as exercises. Using the general property $\sum_{x \in X} (a(x) + b(x)) = \sum_{x \in X} a(x) + \sum_{x \in X} b(x)$ (see Exercise 58), we compute

$$\begin{aligned}\langle f + g, h \rangle &= \sum_{x \in X} (f + g)(x)\bar{h}(x) = \sum_{x \in X} (f(x) + g(x))\bar{h}(x) = \sum_{x \in X} [f(x)\bar{h}(x) + g(x)\bar{h}(x)] \\ &= \sum_{x \in X} f(x)\bar{h}(x) + \sum_{x \in X} g(x)\bar{h}(x) = \langle f, h \rangle + \langle g, h \rangle.\end{aligned}$$

Completeness of the Metric Space $\ell_2(X)$. Let $(f_0, f_1, f_2, \dots) = (f_n : n \in \mathbb{N})$ be a Cauchy sequence in $\ell_2(X)$; we must prove that this sequence converges to some $g \in \ell_2(X)$. To find g , fix $x_0 \in X$ and consider the sequence of complex numbers $(f_n(x_0) : n \in \mathbb{N})$. Given $\epsilon > 0$, choose $m_0 \in \mathbb{N}$ so that for all $m, n \geq m_0$, $d(f_m, f_n) = \sqrt{\sum_{x \in X} |f_m(x) - f_n(x)|^2} < \epsilon$. Then for $m, n \geq m_0$, $|f_m(x_0) - f_n(x_0)| = \sqrt{|f_m(x_0) - f_n(x_0)|^2} \leq d(f_m, f_n) < \epsilon$. So $(f_n(x_0) : n \in \mathbb{N})$ is a Cauchy sequence in the complete metric space \mathbb{C} . Therefore, there is a unique complex number $g(x_0)$ such that $\lim_{n \rightarrow \infty} f_n(x_0) = g(x_0)$. This holds for each $x_0 \in X$, so we have a function $g : X \rightarrow \mathbb{C}$ such that f_n converges to g pointwise. It remains to show that $g \in \ell_2(X)$ and $\lim_{n \rightarrow \infty} f_n = g$ in the metric space $\ell_2(X)$.

To see that $g \in \ell_2(X)$, fix $m_0 \in \mathbb{N}$ (corresponding to $\epsilon = 1$) so that $m, n \geq m_0$ implies $d(f_m, f_n) < 1$. We will show that $(\|f_{m_0}\| + 2)^2$ is a finite upper bound for all sums $\sum_{x \in X'} |g(x)|^2$ as X' ranges over $\mathcal{F}(X)$, so that $\sum_{x \in X} |g(x)|^2 \leq (\|f_{m_0}\| + 2)^2 < \infty$. Fix $X' = \{x_1, \dots, x_N\} \in \mathcal{F}(X)$. Since $f_n(x_k) \rightarrow g(x_k)$ for $k = 1, 2, \dots, N$, we can choose m_1, \dots, m_N such that for all $n \geq m_k$, $|f_n(x_k) - g(x_k)| < 1/\sqrt{N}$. Choose an n larger than all of m_0, m_1, \dots, m_N . Writing $g(x_k) = (g(x_k) - f_n(x_k)) + (f_n(x_k) - f_{m_0}(x_k)) + f_{m_0}(x_k)$ and using (14.3) (extended to a sum of three terms), compute

$$\begin{aligned}\sum_{x \in X'} |g(x)|^2 &\leq \left(\sqrt{\sum_{k=1}^N |g(x_k) - f_n(x_k)|^2} + \sqrt{\sum_{k=1}^N |f_n(x_k) - f_{m_0}(x_k)|^2} + \sqrt{\sum_{k=1}^N |f_{m_0}(x_k)|^2} \right)^2 \\ &\leq \left(\sqrt{\sum_{k=1}^N \frac{1}{N}} + \|f_n - f_{m_0}\| + \|f_{m_0}\| \right)^2 \leq (\|f_{m_0}\| + 2)^2.\end{aligned}$$

The proof that $f_n \rightarrow g$ in $\ell_2(X)$ requires a similar computation. Fix $\epsilon > 0$, and choose $m_0 \in \mathbb{N}$ so that $m, n \geq m_0$ implies $d(f_m, f_n) < \epsilon/2$. We fix $n \geq m_0$ and show that $\|f_n - g\| \leq \epsilon$. To do so, pick $X' = \{x_1, \dots, x_N\}$ in $\mathcal{F}(X)$, and choose m_1, \dots, m_N so that $m \geq m_k$ implies $|f_m(x_k) - g(x_k)| < \epsilon/2\sqrt{N}$. Choose an m larger than all of m_0, m_1, \dots, m_N .

Writing $f_n(x_k) - g(x_k) = (f_n(x_k) - f_m(x_k)) + (f_m(x_k) - g(x_k))$, we compute (using (14.3))

$$\begin{aligned} \sum_{x \in X'} |f_n(x) - g(x)|^2 &\leq \left(\sqrt{\sum_{k=1}^N |f_n(x_k) - f_m(x_k)|^2} + \sqrt{\sum_{k=1}^N |f_m(x_k) - g(x_k)|^2} \right)^2 \\ &< \left(\|f_n - f_m\| + \sqrt{\sum_{k=1}^N \frac{\epsilon^2}{4N}} \right)^2 < (\epsilon/2 + \epsilon/2)^2 = \epsilon^2. \end{aligned}$$

The upper bound of ϵ^2 holds for all X' , so $\sum_{x \in X} |f_n(x) - g(x)|^2 \leq \epsilon^2$, and hence $\|f_n - g\| \leq \epsilon$ as needed.

14.11 Basic Properties of Hilbert Spaces

In this section, we derive some basic algebraic and analytic properties of Hilbert spaces that do not involve the axiom of completeness. Let H be any Hilbert space. For all $x, y \in H$, we have already derived the *triangle inequality for norms*: $\|x+y\| \leq \|x\| + \|y\|$. We can obtain a sharper result called the *Pythagorean theorem* when x and y are orthogonal vectors. This theorem says that if $\langle x, y \rangle = 0$, then $\|x+y\|^2 = \|x\|^2 + \|y\|^2$. Geometrically, the square of the length of the hypotenuse of a right triangle equals the sum of the squares of the two legs. To prove the Pythagorean theorem, we compute (cf. (14.1))

$$\|x+y\|^2 = \langle x+y, x+y \rangle = \langle x, x \rangle + \langle x, y \rangle + \langle y, x \rangle + \langle y, y \rangle = \|x\|^2 + \|y\|^2,$$

where $\langle x, y \rangle = 0$ by hypothesis, and $\langle y, x \rangle = \overline{\langle x, y \rangle} = 0$. More generally, let us say that a list x_1, \dots, x_n of vectors in H is *orthogonal* iff $\langle x_i, x_j \rangle = 0$ for all $i \neq j$ in $[n]$. By induction on $n \in \mathbb{N}^+$, we see that for any orthogonal list x_1, \dots, x_n in H , $\|x_1 + x_2 + \dots + x_n\|^2 = \|x_1\|^2 + \|x_2\|^2 + \dots + \|x_n\|^2$. More generally, if x_1, \dots, x_n is an orthogonal list and $c_1, \dots, c_n \in \mathbb{C}$ are any scalars, then

$$\|c_1x_1 + c_2x_2 + \dots + c_nx_n\|^2 = |c_1|^2\|x_1\|^2 + |c_2|^2\|x_2\|^2 + \dots + |c_n|^2\|x_n\|^2.$$

The key to the induction proof is that c_nx_n is orthogonal to any linear combination $c_1x_1 + \dots + c_{n-1}x_{n-1}$, so that the Pythagorean theorem for the sum of two vectors can be applied.

Here are two more identities resembling the Pythagorean theorem. First, the *parallelogram law* states that for any x, y in a Hilbert space H ,

$$\|x+y\|^2 + \|x-y\|^2 = 2\|x\|^2 + 2\|y\|^2. \quad (14.4)$$

Geometrically, the sum of the squares of the lengths of the two diagonals of any parallelogram equals the sum of the squares of the lengths of the four sides of the parallelogram. To prove this, compute

$$\begin{aligned} \|x+y\|^2 + \|x-y\|^2 &= \langle x+y, x+y \rangle + \langle x-y, x-y \rangle \\ &= [\langle x, x \rangle + \langle x, y \rangle + \langle y, x \rangle + \langle y, y \rangle] + [\langle x, x \rangle - \langle x, y \rangle - \langle y, x \rangle + \langle y, y \rangle] \\ &= 2\|x\|^2 + 2\|y\|^2. \end{aligned}$$

Second, the *polarization identity* states that for all x, y in a Hilbert space H ,

$$\|x + y\|^2 + i\|x + iy\|^2 - \|x - y\|^2 - i\|x - iy\|^2 = 4\langle x, y \rangle. \quad (14.5)$$

This identity is proved by a calculation similar to the one just given (Exercise 60). One can use the polarization identity to prove that any complex normed vector space (as defined in §10.4) that is complete and whose norm satisfies the parallelogram law must be a Hilbert space (Exercise 61).

Next we discuss some continuity properties of the operations appearing in the definition of a Hilbert space. First, *addition is continuous*: if $x_n \rightarrow x$ and $y_n \rightarrow y$ in a Hilbert space H , then $x_n + y_n \rightarrow x + y$. To prove this, fix $\epsilon > 0$ and choose $n_0, n_1 \in \mathbb{N}$ such that $n \geq n_0$ implies $\|x_n - x\| < \epsilon/2$, whereas $n \geq n_1$ implies $\|y_n - y\| < \epsilon/2$. Then for $n \geq \max(n_0, n_1)$,

$$\|(x_n + y_n) - (x + y)\| = \|(x_n - x) + (y_n - y)\| \leq \|x_n - x\| + \|y_n - y\| < \epsilon.$$

Second, *scalar multiplication is continuous*: if $x_n \rightarrow x$ in H and $c_n \rightarrow c$ in \mathbb{C} , then $c_n x_n \rightarrow cx$ in H . The convergent sequence (c_n) must be bounded in \mathbb{C} ; say $|c_n| \leq M$ for all $n \in \mathbb{N}$ for some constant $M > 0$. Now, given $\epsilon > 0$, choose $n_0, n_1 \in \mathbb{N}$ so that $n \geq n_0$ implies $\|x_n - x\| < \epsilon/(2M)$, whereas $n \geq n_1$ implies $|c_n - c| < \epsilon/(2(1 + \|x\|))$. Then $n \geq \max(n_0, n_1)$ implies

$$\|c_n x_n - cx\| = |c_n(x_n - x) + (c_n - c)x| \leq |c_n| \cdot \|x_n - x\| + |c_n - c| \cdot \|x\| < M(\epsilon/(2M)) + \epsilon/2 = \epsilon.$$

Third, *the inner product on H is continuous*: if $x_n \rightarrow x$ and $y_n \rightarrow y$ in H , then $\langle x_n, y_n \rangle \rightarrow \langle x, y \rangle$ in \mathbb{C} . As before, there is a single constant $M > 0$ such that $\|x_n\|$, $\|x\|$, $\|y_n\|$, and $\|y\|$ are all bounded above by M . Given $\epsilon > 0$, choose $n_0, n_1 \in \mathbb{N}$ so that $n \geq n_0$ implies $\|x_n - x\| < \epsilon/(2M)$, whereas $n \geq n_1$ implies $\|y_n - y\| < \epsilon/(2M)$. Then $n \geq \max(n_0, n_1)$ implies

$$|\langle x_n, y_n \rangle - \langle x, y \rangle| = |\langle x_n - x, y_n \rangle + \langle x, y_n - y \rangle| \leq |\langle x_n - x, y_n \rangle| + |\langle x, y_n - y \rangle|.$$

By the Cauchy–Schwarz inequality in H ,

$$|\langle x_n - x, y_n \rangle| + |\langle x, y_n - y \rangle| \leq \|x_n - x\| \cdot \|y_n\| + \|x\| \cdot \|y_n - y\| < (\epsilon/(2M))M + M(\epsilon/(2M)) = \epsilon.$$

Fourth, *the norm on H is continuous*: if $x_n \rightarrow x$ in H , then $\|x_n\| \rightarrow \|x\|$ in \mathbb{R} . To prove this, note $\|x_n\| = \|x_n - x + x\| \leq \|x_n - x\| + \|x\|$ implies $\|x_n\| - \|x\| \leq \|x_n - x\|$. Similarly, $\|x\| - \|x_n\| \leq \|x - x_n\| = \|x_n - x\|$, so the absolute value of $\|x_n\| - \|x\|$ is at most $\|x_n - x\|$. Given $\epsilon > 0$, choose n_0 so $n \geq n_0$ implies $\|x_n - x\| < \epsilon$. Then $n \geq n_0$ implies $|\|x_n\| - \|x\|| < \epsilon$, so that $\|x_n\| \rightarrow \|x\|$ in \mathbb{R} .

14.12 Closed Convex Sets in Hilbert Spaces

Recall from Chapter 11 that a subset C of a (real or complex) vector space V is called *convex* iff for all $x, y \in C$ and all real $t \in [0, 1]$, $tx + (1 - t)y \in C$. Geometrically, this condition says that a convex set must contain the line segment joining any two of its points. Every subspace W of V is a convex set. If C is a convex subset of V and $z \in V$, the translate $z + C = \{z + w : w \in C\}$ is convex. The empty set is also convex.

The following geometric lemma will turn out to be a key technical tool for studying Hilbert spaces. *For every nonempty, closed, convex set C in a Hilbert space H , there exists*

a unique $x \in C$ of minimum norm; i.e., $\|x\| < \|z\|$ for all $z \neq x$ in C . We prove uniqueness first. Suppose $x, y \in C$ are two elements such that $r = \|x\| = \|y\| \leq \|z\|$ for all $z \in C$; we must prove $x = y$. Consider $z = (1/2)x + (1/2)y$. On one hand, by convexity of C , $z \in C$ and hence $r \leq \|z\|$. On the other hand, applying the parallelogram law to $x/2$ and $y/2$ shows that

$$\|(x/2) + (y/2)\|^2 + \|(x/2) - (y/2)\|^2 = 2\|(x/2)\|^2 + 2\|(y/2)\|^2, \quad (14.6)$$

which simplifies to $\|z\|^2 + (1/4)\|x - y\|^2 = (2/4)r^2 + (2/4)r^2 = r^2$. Then $\|x - y\|^2 = 4r^2 - 4\|z\|^2 \leq 4r^2 - 4r^2 = 0$, forcing $\|x - y\| = 0$ and $x = y$.

Turning to the existence proof, let s be the greatest lower bound in \mathbb{R} of the set $\{\|z\| : z \in C\}$; this set is nonempty and bounded below by zero, so s does exist. By definition of greatest lower bound, for each $n \in \mathbb{N}^+$ we can find $x_n \in C$ such that $s \leq \|x_n\| < s + 1/n$. Observe that $\lim_{n \rightarrow \infty} \|x_n\| = s$. We first show that $(x_n : n \in \mathbb{N}^+)$ is a Cauchy sequence in H . Fix $\epsilon > 0$; we must find $m_0 \in \mathbb{N}^+$ so that $m, n \geq m_0$ implies $d(x_n, x_m) = \|x_n - x_m\| < \epsilon$. For any $m, n \in \mathbb{N}^+$, note that $(1/2)x_m + (1/2)x_n \in C$ by convexity of C , so that $s \leq \|(x_m/2) + (x_n/2)\|$ by definition of s . Applying the parallelogram law to $x_n/2$ and $x_m/2$ gives

$$\begin{aligned} \|(x_n/2) - (x_m/2)\|^2 &= 2\|x_n/2\|^2 + 2\|x_m/2\|^2 - \|(x_m/2) + (x_n/2)\|^2 \\ &\leq \frac{(s + 1/n)^2}{2} + \frac{(s + 1/m)^2}{2} - s^2, \end{aligned}$$

which rearranges to

$$\|x_n - x_m\| \leq 2\sqrt{s/n + s/m + 1/(2n^2) + 1/(2m^2)}.$$

Since s is fixed, each term inside the square root approaches zero as n and m increase to infinity. So we can choose m_0 so that $m, n \geq m_0$ implies $\|x_n - x_m\| < \epsilon$.

Now we know $(x_n : n \in \mathbb{N}^+)$ is a Cauchy sequence. By completeness of the Hilbert space H , this sequence must converge to some point $y \in H$. Because C is a closed subset of H , we must have $y \in C$. By continuity of the norm, $s = \lim_{n \rightarrow \infty} \|x_n\| = \|\lim_{n \rightarrow \infty} x_n\| = \|y\|$. So $\|y\|$ is a lower bound of all the norms $\|x\|$ for $x \in C$, completing the existence proof.

Here is a slight generalization of the lemma just proved. *For every nonempty, closed, convex set C in H and all $w \in H$, there exists a unique $x \in C$ minimizing $d(x, w)$* ; we call x the point of C closest to w . To prove this, consider the translate $C' = (-w) + C$, which is readily verified to be nonempty, closed, and convex (Exercise 64). The map $x \mapsto x - w$ is a bijection between C and C' . Moreover, $d(x, w) = \|x - w\|$. Thus, $x \in C$ will minimize $d(x, w)$ iff $x - w \in C'$ has minimum norm among all elements of C' . By the lemma already proved, there exists a unique element of C' with the latter property. Therefore, there exists a unique $x \in C$ minimizing $d(x, w)$, as needed.

14.13 Orthogonal Complements

In our study of *finite-dimensional* inner product spaces in Chapter 13, we introduced the idea of the *orthogonal complement* of a subspace W , denoted W^\perp . For each subspace W of the inner product space V , W^\perp is the subspace consisting of all vectors in V that are orthogonal to every vector in W . We showed that the map $W \mapsto W^\perp$ is an inclusion-reversing bijection on the lattice of subspaces of V , and for every subspace W , $W^{\perp\perp} = W$.

and $V = W \oplus W^\perp$. A key ingredient in proving these results was the dimension formula $\dim(V) = \dim(W) + \dim(W^\perp)$, whose proof required V to be finite-dimensional.

We wish to extend these results to the case of a general Hilbert space H , which may be infinite-dimensional. To begin, we define the *orthogonal complement* of an arbitrary subset S of H to be

$$S^\perp = \{v \in H : \langle v, w \rangle = 0 \text{ for all } w \in S\}.$$

We claim that S^\perp is always a *closed subspace* of H . First, $0_H \in S^\perp$ since $\langle 0, w \rangle = 0$ for all $w \in S$. Second, given $u, v \in S^\perp$, we know $\langle u, w \rangle = 0 = \langle v, w \rangle$ for all $w \in S$. So $\langle u+v, w \rangle = \langle u, w \rangle + \langle v, w \rangle = 0+0=0$ for all $w \in S$, and $u+v \in S^\perp$. Third, given $u \in S^\perp$ and $c \in \mathbb{C}$, we find that $\langle cu, w \rangle = c\langle u, w \rangle = c0 = 0$ for all $w \in S$, so $cu \in S^\perp$.

Fourth, to see that S^\perp is closed, define a map $R_w : H \rightarrow \mathbb{C}$ (for each $w \in H$) by letting $R_w(x) = \langle x, w \rangle$ for all $x \in H$. If $x_n \rightarrow x$ in H , then we saw in §14.11 that $R_w(x_n) = \langle x_n, w \rangle \rightarrow \langle x, w \rangle = R_w(x)$. So, R_w is a *continuous* map from H to \mathbb{C} . In particular, since the one-point set $\{0\}$ is closed in \mathbb{C} , the inverse image $R_w^{-1}[\{0\}] = \{x \in H : \langle x, w \rangle = 0\}$ is a closed subset of H . By definition of S^\perp , we have $S^\perp = \bigcap_{w \in S} R_w^{-1}[\{0\}]$. This is an intersection of a family of closed sets, so S^\perp is closed as claimed. We remark that each R_w is a \mathbb{C} -linear map (by the inner product axioms), and $R_w^{-1}[\{0\}]$ is precisely the kernel of R_w . This gives another way to see that S^\perp is a subspace, since the kernel of a linear map is a subspace of the domain, and the intersection of a family of subspaces is also a subspace.

To obtain a bijective correspondence $W \mapsto W^\perp$ in the setting of Hilbert spaces, we must restrict attention to the set of *closed* subspaces W , since applying the orthogonal complement operator always produces a closed subspace. We intend to show that *for any closed subspace W of a Hilbert space H , $H = W \oplus W^\perp$ and $W^{\perp\perp} = W$* . Dimension-counting arguments are no longer available, but we can instead invoke the geometric lemma of §14.12, whose proof made critical use of the completeness of H .

Given $x \in H$, we must prove there exist unique $y \in W$ and $z \in W^\perp$ with $x = y + z$. We prove uniqueness first: assume $y_1, y_2 \in W$ and $z_1, z_2 \in W^\perp$ satisfy $x = y_1 + z_1 = y_2 + z_2$. Let $u = y_1 - y_2 = z_2 - z_1$. Since W is a subspace, $u = y_1 - y_2 \in W$. Since W^\perp is a subspace, $u = z_2 - z_1 \in W^\perp$. Then $\langle u, u \rangle = 0$, forcing $u = 0$ by the inner product axioms. So $y_1 = y_2$ and $z_1 = z_2$.

To prove existence of y and z , we can draw intuition from the case where W is a two-dimensional subspace of \mathbb{R}^3 and W^\perp is the line through 0 perpendicular to W . In this case, we could find y given x by dropping an altitude from x to the plane W . This altitude meets W at the point y on that plane closest to x , and then one has $x = y + (x - y)$ where the vector $x - y$ is parallel to the altitude and hence is in W^\perp . This suggests that in the general case, we could define y to be the unique point in W closest to x , and let $z = x - y$. The point y does exist, since W is a nonempty closed convex subset of H . It is evident that $x = y + z$, but we must still check that $z \in W^\perp$.

Fix $w \in W$; we must show $\langle z, w \rangle = 0$. The conclusion holds for $w = 0$, so assume $w \neq 0$. Write $w = cw$, where $c = \|w\|$ and $u = c^{-1}w \in W$ satisfies $\|u\| = 1$. For any $s \in \mathbb{C}$, the vector $y - su$ lies in the subspace W . Since y is the closest point in W to x , we have $\|z\| = \|x - y\| \leq \|x - (y - su)\| = \|z + su\|$ for all $s \in \mathbb{C}$. Squaring this inequality and rewriting using scalar products,

$$\langle z, z \rangle \leq \langle z, z \rangle + \bar{s}\langle z, u \rangle + s\langle u, z \rangle + |s|^2\langle u, u \rangle.$$

Since $\langle u, u \rangle = 1$ and $\langle u, z \rangle = \overline{\langle z, u \rangle}$, the inequality becomes $0 \leq \bar{s}\langle z, u \rangle + s\overline{\langle z, u \rangle} + |s|^2$ for all $s \in \mathbb{C}$. Choose $s = -\langle z, u \rangle$ to get

$$0 \leq -|\langle z, u \rangle|^2 - |\langle z, u \rangle|^2 + |\langle z, u \rangle|^2 = -|\langle z, u \rangle|^2,$$

which forces $\langle z, u \rangle = 0$. Then $\langle z, w \rangle = c\langle z, u \rangle = 0$, as needed.

We have now proved $H = W \oplus W^\perp$ for any closed subspace W ; we use this to prove $W^{\perp\perp} = W$. Recalling that $\langle x, y \rangle = 0$ iff $\langle y, x \rangle = 0$, we see from the definitions that $W \subseteq W^{\perp\perp}$ without any hypothesis on the subset W . The result just proved shows that $H = W \oplus W^\perp$ and also $H = W^\perp \oplus W^{\perp\perp}$, since W^\perp is a closed subspace. We use this to prove $W^{\perp\perp} \subseteq W$, as follows. Fix $x \in W^{\perp\perp}$. There exist unique $y \in W$ and $z \in W^\perp$ with $x = y + z$. Similarly, x can be written in exactly one way as the sum of a vector in $W^{\perp\perp}$ and a vector in W^\perp . Since $y \in W^{\perp\perp}$, one such sum is $x = y + z$. On the other hand, since $0 \in W^\perp$, another such sum is $x = x + 0$. By uniqueness, this forces $y = x$ and $z = 0$, so $x = y$ is in W , as needed.

Let L be the set of all closed subspaces of the given Hilbert space H ; one checks that L is a complete lattice ordered by set inclusion. We have shown that $f : L \rightarrow L$, given by $f(W) = W^\perp$ for $W \in L$, does map into the codomain L and satisfies $f(f(W)) = W$. Therefore, f is a bijection on L with $f^{-1} = f$. One readily verifies that for $W, X \in L$, $W \subseteq X$ implies $f(W) \supseteq f(X)$, so f is order-reversing. Thus we have shown that $W \mapsto W^\perp$ is a lattice anti-isomorphism of the lattice of all closed subspaces of H .

14.14 Orthonormal Sets

Orthonormal bases play a central role in finite-dimensional inner product spaces. Every such space has an orthonormal basis, and every vector in the space can be written as a (finite) linear combination of these basis elements. In the setting of Hilbert spaces, we will develop an analytic version of orthonormal bases in which “infinite linear combinations” of basis elements are allowed. To prepare for this, we first study finite orthonormal sets in a general Hilbert space H .

A subset X of H is called *orthonormal* iff $\|x\|^2 = \langle x, x \rangle = 1$ for all $x \in X$, and $\langle x, y \rangle = 0$ for all $x \neq y$ in X . An orthonormal set X is automatically linearly independent. For suppose $\{x_1, \dots, x_N\}$ is any finite subset of X and $c_1, \dots, c_N \in \mathbb{C}$ satisfy $c_1x_1 + \dots + c_Nx_N = 0$. Then for $1 \leq j \leq N$,

$$0 = \langle 0, x_j \rangle = \langle c_1x_1 + \dots + c_Nx_N, x_j \rangle = c_j\langle x_j, x_j \rangle + \sum_{k \neq j} c_k\langle x_k, x_j \rangle = c_j.$$

Now suppose $X = \{x_1, \dots, x_N\}$ is a finite orthonormal subset of H . Let W be the subspace of H spanned by X . We will give a direct argument to show that $H = W \oplus W^\perp$. (This also follows from previous results and the fact that the finite-dimensional subspace W must be closed.) It is routine to check that $W \cap W^\perp = \{0\}$, since 0 is the only vector orthogonal to itself. Next, we show how to write any $x \in H$ in the form $x = y + z$, where $y \in W$ and $z \in W^\perp$. Define $y = \sum_{k=1}^N \langle x, x_k \rangle x_k \in W$ and $z = x - y$. To prove $z \in W^\perp$, it suffices (by the inner product axioms) to show that $\langle z, x_j \rangle = 0$ for $1 \leq j \leq N$. We compute

$$\langle z, x_j \rangle = \langle x, x_j \rangle - \langle y, x_j \rangle = \langle x, x_j \rangle - \sum_{k=1}^N \langle x, x_k \rangle \langle x_k, x_j \rangle = \langle x, x_j \rangle - \langle x, x_j \rangle = 0.$$

Thus $H = W \oplus W^\perp$. We already know $H = W^{\perp\perp} \oplus W^\perp$, so the argument used at the end of §14.13 shows that $W = W^{\perp\perp}$, and hence W is a closed subspace. It now follows from the proof in §14.13 that $y = \sum_{k=1}^N \langle x, x_k \rangle x_k$ is the *closest* element of W to x . We call y the *orthogonal projection* of x onto W .

Since y is orthogonal to z and the x_k 's are orthogonal to each other, the Pythagorean theorem shows that $\|x\|^2 = \|y\|^2 + \|z\|^2 = \sum_{k=1}^N |\langle x, x_k \rangle|^2 + \|z\|^2$. Discarding the “error term” $\|z\|^2$, we obtain the inequality

$$\sum_{k=1}^N |\langle x, x_k \rangle|^2 \leq \|x\|^2,$$

which is the finite version of *Bessel's inequality*. This inequality can be viewed as an approximate version of the Pythagorean theorem: the sum of the squared norms of the components of x in certain orthogonal directions is at most the squared length of x itself.

Next let X be any orthonormal set (possibly infinite) in a Hilbert space H . Given $w \in H$, define a function $f_w : X \rightarrow \mathbb{C}$ by setting $f_w(x) = \langle w, x \rangle$ for all $x \in X$. The complex scalars $\langle w, x \rangle$ are called the *Fourier coefficients* of w relative to the orthonormal set X . We claim that for all $w \in H$, f_w lies in the space $\ell_2(X)$ of square-summable sequences indexed by X (see §14.9). In other words, $\sum_{x \in X} |f_w(x)|^2 = \sum_{x \in X} |\langle w, x \rangle|^2 < \infty$.

To verify the claim, fix any finite subset $X' = \{x_1, \dots, x_N\}$ of X , which is also orthonormal. By the finite version of Bessel's inequality,

$$\sum_{x \in X'} |f_w(x)|^2 = \sum_{k=1}^N |\langle w, x_k \rangle|^2 \leq \|w\|^2.$$

Thus $\|w\|^2$ is a finite upper bound for all these sums as X' ranges over $\mathcal{F}(X)$. Thus we obtain the general version of Bessel's inequality, namely

$$\sum_{x \in X} |\langle w, x \rangle|^2 \leq \|w\|^2$$

for all $w \in H$ and all orthonormal sets $X \subseteq H$.

14.15 Maximal Orthonormal Sets

A *maximal orthonormal set* in a Hilbert space H is an orthonormal set X such that for any set Y properly containing X , Y is not orthonormal. Some texts refer to maximal orthonormal sets as *complete* orthonormal sets; but we avoid this term to prevent confusion with the notion of a complete metric space. By appealing to Zorn's Lemma (Exercise 69), one sees that *maximal orthonormal sets exist in any Hilbert space*.

Let X be a maximal orthonormal set in a Hilbert space H . We will show that in this case, equality holds in Bessel's inequality, i.e.,

$$\sum_{x \in X} |\langle w, x \rangle|^2 = \|w\|^2 \tag{14.7}$$

for all $w \in H$. To get a contradiction, assume this equality fails for some $w \in H$. Let $r = \sum_{x \in X} |\langle w, x \rangle|^2 < \|w\|^2$. For each $n \in \mathbb{N}^+$, we can find a finite subset X_n of X such that $r - 1/n < \sum_{x \in X_n} |\langle w, x \rangle|^2 \leq r$. By replacing each X_n by $X_1 \cup X_2 \cup \dots \cup X_n$ (which is still finite), we can arrange that $X_1 \subseteq X_2 \subseteq \dots \subseteq X_n \subseteq \dots$. Now define x_n and y_n in H by setting $x_n = \sum_{x \in X_n} \langle w, x \rangle x$ and $y_n = w - x_n$ for $n \in \mathbb{N}^+$.

We claim (x_n) is a Cauchy sequence in H . Given $\epsilon > 0$, choose $m_0 \in \mathbb{N}^+$ so that

$1/m_0 < \epsilon$. For $m \geq n \geq m_0$, the Pythagorean theorem gives

$$\|x_m - x_n\|^2 = \left\| \sum_{x \in X_m \sim X_n} \langle w, x \rangle x \right\|^2 = \sum_{x \in X_m} |\langle w, x \rangle|^2 - \sum_{x \in X_n} |\langle w, x \rangle|^2 < r - (r - 1/n) < \epsilon.$$

So (x_n) is Cauchy, hence converges to a point $z \in H$ by completeness of H . Letting $y = w - z$, we have $y_n = (w - x_n) \rightarrow (w - z) = y$ as n goes to infinity.

Now, $\|z\| = \lim_{n \rightarrow \infty} \|x_n\| = \lim_{n \rightarrow \infty} \sqrt{\sum_{x \in X_n} |\langle w, x \rangle|^2} = \sqrt{r} < \|w\|$. It follows that $z \neq w$ and $y \neq 0$. We next claim that $\langle y, x \rangle = 0$ for all $x \in X$. Once this claim is proved, we can deduce that $y \notin X$, $y/\|y\| \notin X$, yet $X \cup \{y/\|y\|\}$ is orthonormal, contradicting the maximality of the orthonormal set X . To prove the claim, fix $x \in X$ and note that $\langle y, x \rangle = \lim_{n \rightarrow \infty} \langle y_n, x \rangle = \langle w, x \rangle - \lim_{n \rightarrow \infty} \langle x_n, x \rangle$. If $x \in X_{n_0}$ for some n_0 , then for all $n \geq n_0$, x_n is the sum of $\langle w, x \rangle x$ plus other vectors orthogonal to x , so $\lim_{n \rightarrow \infty} \langle x_n, x \rangle = \langle w, x \rangle$ and $\langle y, x \rangle = 0$. On the other hand, if $x \notin X_{n_0}$ for all n_0 , then $\langle x_n, x \rangle = 0$ for all n . We will show that $\langle w, x \rangle = 0$ in this case. If $\langle w, x \rangle \neq 0$, choose n so that $|\langle w, x \rangle|^2 > 1/n$. Then $|\langle w, x \rangle|^2 + \sum_{u \in X_n} |\langle w, u \rangle|^2 > 1/n + (r - 1/n) = r$, which contradicts the definition of r . We have now proved that $\|f_w\|^2 = \sum_{x \in X} |\langle w, x \rangle|^2 = \|w\|^2$ for all $w \in H$.

We are now ready to introduce the idea of an “infinite linear combination” of vectors in an orthonormal set. We will show that for any maximal orthonormal set X and all $w \in H$, $w = \sum_{x \in X} \langle w, x \rangle x$. By definition, this means that given any $\epsilon > 0$, there exists a finite subset $X' \in \mathcal{F}(X)$ such that for every finite subset $Y \in \mathcal{F}(X)$ containing X' , $\|w - \sum_{x \in Y} \langle w, x \rangle x\| < \epsilon$. Given $w \in H$ and $\epsilon > 0$, define the sets X_n as in the proof above (taking $r = \|w\|^2$ here). Choose n with $1/n < \epsilon^2$, and take X' to be the finite set X_n . For any finite subset Y containing X_n , we know from the calculation in §14.14 that

$$\left\| w - \sum_{x \in Y} \langle w, x \rangle x \right\|^2 = \|w\|^2 - \sum_{x \in Y} |\langle w, x \rangle|^2 \leq \|w\|^2 - \sum_{x \in X_n} |\langle w, x \rangle|^2 < 1/n < \epsilon^2,$$

as needed.

14.16 Isomorphism of H and $\ell_2(X)$

Recall that a *vector space isomorphism* is a bijection $f : V \rightarrow W$ between F -vector spaces V and W such that $f(x+y) = f(x) + f(y)$ and $f(cx) = cf(x)$ for all $x, y \in V$ and all $c \in F$. Given two metric spaces X and Y , an *isometry* from X to Y is a function $f : X \rightarrow Y$ that “preserves distances,” i.e., $d_Y(f(u), f(v)) = d_X(u, v)$ for all $u, v \in X$. An isometry must be one-to-one, since if $u, v \in X$ satisfy $f(u) = f(v)$, then $d_X(u, v) = d_Y(f(u), f(v)) = 0$ and hence $u = v$. Similarly, one checks that *an isometry must be continuous* (Exercise 30). Given Hilbert spaces H_1 and H_2 , a *Hilbert space isomorphism* (also called an *isometric isomorphism*) is a bijection $f : H_1 \rightarrow H_2$ that is both a vector space isomorphism and an isometry.

A Hilbert space isomorphism preserves norms, i.e., $\|f(u)\| = \|u\|$ for all $u \in H$. For, $\|f(u)\| = d(f(u), 0) = d(f(u), f(0)) = d(u, 0) = \|u\|$. It now follows from the polarization identity that a Hilbert space isomorphism will automatically preserve inner products, i.e.,

$\langle f(u), f(v) \rangle = \langle u, v \rangle$ for all $u, v \in H_1$. To see why, use (14.5) in H_2 and in H_1 to compute

$$\begin{aligned}\langle f(u), f(v) \rangle &= \frac{1}{4}(\|f(u) + f(v)\|^2 - \|f(u) - f(v)\|^2 + i\|f(u) + if(v)\|^2 - i\|f(u) - if(v)\|^2) \\ &= \frac{1}{4}(\|f(u+v)\|^2 - \|f(u-v)\|^2 + i\|f(u+iv)\|^2 - i\|f(u-iv)\|^2) \\ &= \frac{1}{4}(\|u+v\|^2 - \|u-v\|^2 + i\|u+iv\|^2 - i\|u-iv\|^2) = \langle u, v \rangle.\end{aligned}$$

Conversely, if a given linear map $f : H_1 \rightarrow H_2$ preserves inner products, one sees by taking $u = v$ that f preserves norms. By linearity, f must be an isometry.

Let H be any Hilbert space, and let X be any maximal orthonormal subset of H . Such subsets do exist, by Exercise 69. We define a map $f : H \rightarrow \ell_2(X)$ by letting $f(w) = f_w$ for all $w \in H$. Recall that $f_w : X \rightarrow \mathbb{C}$ gives the Fourier coefficients of w relative to X , namely $f_w(x) = \langle w, x \rangle$ for $x \in X$. Also recall that f_w does lie in $\ell_2(X)$, by Bessel's inequality. Our goal is to prove that f is an isometric isomorphism, so that *every Hilbert space H is isomorphic to a Hilbert space of the form $\ell_2(X)$, where X can be chosen to be a maximal orthonormal subset of H .*

First, we check \mathbb{C} -linearity of f . Fix $w, y \in H$ and $c \in \mathbb{C}$. On one hand, $f(w+y) = f_{w+y}$ is the function sending $x \in X$ to $\langle w+y, x \rangle = \langle w, x \rangle + \langle y, x \rangle$. On the other hand, $f(w) + f(y) = f_w + f_y$ is the function sending $x \in X$ to $f_w(x) + f_y(x) = \langle w, x \rangle + \langle y, x \rangle$. These functions are equal, so $f(w+y) = f(w) + f(y)$. Similarly, both functions $f(cw)$ and $cf(w)$ send each $x \in X$ to $\langle cw, x \rangle = c\langle w, x \rangle$, so $f(cw) = cf(w)$. Next, we observe that f is an isometry, since (14.7) says

$$\|w\| = \sqrt{\sum_{x \in X} |\langle w, x \rangle|^2} = \sqrt{\sum_{x \in X} |f_w(x)|^2} = \|f_w\|$$

for all $w \in H$. We deduce that f is one-to-one. It also follows that f preserves inner products, so that

$$\langle f_w, f_z \rangle = \sum_{x \in X} \langle w, x \rangle \overline{\langle z, x \rangle} = \langle w, z \rangle$$

for all $w, z \in H$. This formula is often called *Parseval's identity*.

The only thing left to prove is that f is surjective. Fix $g \in \ell_2(X)$; we must find $w \in H$ with $f(w) = g$. To do so, we first build a sequence of “partial sums” $g_n \in \ell_2(X)$ such that $\lim_{n \rightarrow \infty} g_n = g$. We know $r = \sum_{x \in X} |g(x)|^2 < \infty$. So, for each $n \in \mathbb{N}^+$, there is a finite subset X_n of X such that $r - 1/n < \sum_{x \in X_n} |g(x)|^2 \leq r$. Define $g_n(x) = g(x)$ for $x \in X_n$, and $g_n(x) = 0$ for $x \in X \setminus X_n$. Evidently, $g_n \in \ell_2(X)$, and $\|g - g_n\|^2 = \sum_{x \in X \setminus X_n} |g(x)|^2 < 1/n$. It readily follows that $g_n \rightarrow g$ in the Hilbert space $\ell_2(X)$. In particular, the convergent sequence (g_n) is also a Cauchy sequence.

For each $n \in \mathbb{N}^+$, define $w_n \in H$ by $w_n = \sum_{x \in X_n} g(x)x$. By orthonormality of X , $\langle w_n, x \rangle = g(x) = g_n(x)$ for $x \in X_n$, and $\langle w_n, x \rangle = 0 = g_n(x)$ for $x \in X \setminus X_n$. Therefore, $f(w_n) = g_n$ for all $n \in \mathbb{N}^+$. Since f is an isometry and (g_n) is a Cauchy sequence, (w_n) must also be a Cauchy sequence (Exercise 31). By completeness of H , w_n converges to some point $w \in H$. Now f is continuous (being an isometry), so $w_n \rightarrow w$ in H implies $g_n = f(w_n) \rightarrow f(w)$ in $\ell_2(X)$. On the other hand, by construction, $g_n \rightarrow g$ in $\ell_2(X)$. Since limits are unique, $g = f(w)$ as needed.

It can be shown that for all sets X and Y , the Hilbert spaces $\ell_2(X)$ and $\ell_2(Y)$ are isomorphic iff $|X| = |Y|$, i.e., iff there is a bijection from X onto Y . So, the theorem in this section provides a *classification of Hilbert spaces* analogous to the classification of F -vector spaces as direct sums of copies of F , where the number of direct summands is the dimension of the vector space.

14.17 Continuous Linear Maps

We intend to study operators and linear functionals on Hilbert spaces. Before doing so, we establish some facts about continuous linear maps in the more general setting of normed vector spaces. Recall from Chapter 10 that a *normed vector space* consists of a real or complex vector space V and a norm function $\|\cdot\| : V \rightarrow \mathbb{R}$ satisfying these axioms: for all $x \in V$, $0 \leq \|x\| < \infty$; for all $x \in V$, if $\|x\| = 0$, then $x = 0$; for all $x, y \in V$, $\|x + y\| \leq \|x\| + \|y\|$ (the triangle inequality); and for all $x \in V$ and all scalars c , $\|cx\| = |c| \cdot \|x\|$. A normed vector space becomes a metric space with distance function $d(x, y) = \|x - y\|$ for $x, y \in V$. A *Banach space* is a normed vector space that is complete relative to this metric.

Consider a map $T : V \rightarrow W$ between two normed vector spaces. Recall that T is *continuous* iff whenever $v_n \rightarrow v$ in V , $T(v_n) \rightarrow T(v)$ in W . To check continuity of a *linear* map T , it suffices to check that for all sequences (z_n) of points in V such that $z_n \rightarrow 0_V$ in V , $T(z_n) \rightarrow 0_W$ in W . For suppose (v_n) is a sequence in V converging to $v \in V$. Then the sequence $(v_n - v)$ converges to $v - v = 0$, so the stated condition on T implies that $T(v_n - v)$ converges to 0_W . By linearity, $T(v_n) - T(v)$ converges to 0_W , so $T(v_n) \rightarrow T(v)$, and T is continuous.

A linear map $T : V \rightarrow W$ is called *bounded* iff there is a finite real constant $M \geq 0$ such that $\|T(x)\|_W \leq M\|x\|_V$ for all $x \in V$. Let us show that a *linear map T is continuous iff it is bounded*. On one hand, suppose T is bounded and $x_n \rightarrow 0$ in V . If $M = 0$, then $T(x) = 0_W$ for all $x \in V$, so T is continuous. Otherwise, given $\epsilon > 0$, choose n_0 so that $n \geq n_0$ implies $\|x_n\| < \epsilon/M$. Then for $n \geq n_0$, $\|T(x_n)\|_W \leq M\|x_n\|_V < \epsilon$. So $T(x_n) \rightarrow 0$, and T is continuous. On the other hand, suppose T is not bounded. Then for each integer $n > 0$, there is $x_n \in V$ (necessarily nonzero) with $\|T(x_n)\|_W > n\|x_n\|_V$. By linearity of T and properties of norms, this inequality will still hold if we replace x_n by any nonzero scalar multiple of itself. Picking an appropriate scalar multiple, we can arrange that $\|x_n\|_V = 1/n$ for all $n \in \mathbb{N}^+$. Then $x_n \rightarrow 0$ in V , but $\|T(x_n)\|_W > 1$ for all $n \in \mathbb{N}^+$, so $T(x_n)$ does not converge to zero in W , and T is not continuous.

Given two normed vector spaces V and W , let $B(V, W)$ be the set of all bounded (i.e., continuous) linear maps from V to W . One checks readily that $B(V, W)$ is a subspace of the vector space of all linear maps from V to W under pointwise operations on functions (see §4.2). We now show that $B(V, W)$ is itself a normed vector space using the *operator norm* defined by

$$\|T\| = \sup\{\|T(x)\|_W : x \in V, \|x\|_V = 1\}$$

for $T \in B(V, W)$. (In the special case $V = \{0\}$, define $\|0_{B(V, W)}\| = 0$.) First, since the continuous map T is bounded, we know there is a finite upper bound M for the set on the right side, so that $\|T\|$ is a finite nonnegative real number. Second, if $\|T\| = 0$, then $\|T(x)\|_W = 0$ for all $x \in V$ with $\|x\|_V = 1$. Multiplying by a scalar, it follows that $\|T(y)\|_W = 0$ for all $y \in V$, so $T(y) = 0_W$ for all $y \in V$, so T is the zero map, which is the zero element of $B(V, W)$. The axiom $\|cT\| = |c| \cdot \|T\|$ follows readily from the fact that $\|(cT)(x)\|_W = \|c(T(x))\|_W = |c| \cdot \|T(x)\|_W$. Finally, for $S, T \in B(V, W)$ and any $x \in V$ of norm 1,

$$\|(S + T)(x)\|_W = \|S(x) + T(x)\|_W \leq \|S(x)\|_W + \|T(x)\|_W \leq \|S\| + \|T\|,$$

so that $\|S + T\| \leq \|S\| + \|T\|$. One readily verifies that:

$$\begin{aligned} \|T(y)\|_W &\leq \|T\| \cdot \|y\|_V \text{ for all } y \in V; \\ \|T\| &= \sup\{\|T(x)\|_W / \|x\|_V : x \in V, x \neq 0\}; \\ \|T\| &= \inf\{M \in \mathbb{R}^+ : \|T(y)\|_W \leq M\|y\|_V \text{ for all } y \in V\}; \end{aligned}$$

and if $U \in B(W, Z)$, then $U \circ T \in B(V, Z)$ and $\|U \circ T\| \leq \|U\| \cdot \|T\|$.

We conclude this section by showing that *if W is a Banach space, then $B(V, W)$ is a Banach space*. In other words, completeness of the metric space W implies completeness of $B(V, W)$. Let (T_n) be a Cauchy sequence in $B(V, W)$. For each fixed $x \in V$, $\|T_n(x) - T_m(x)\|_W = \|(T_n - T_m)(x)\|_W \leq \|T_n - T_m\| \cdot \|x\|_V$. It follows that $(T_n(x))$ is a Cauchy sequence in W for each fixed $x \in V$. By completeness of W , for each $x \in V$ there exists a unique $y \in W$ (denoted $T(x)$) so that $\lim_{n \rightarrow \infty} T_n(x) = y = T(x)$.

We now have a function $T : V \rightarrow W$; we must check that T is linear, T is bounded, and $T_n \rightarrow T$ in $B(V, W)$. For linearity, fix $x, z \in V$ and a scalar c . By continuity and linearity of each T_n ,

$$T(x + z) = \lim_n T_n(x + z) = \lim_n [T_n(x) + T_n(z)] = \lim_n T_n(x) + \lim_n T_n(z) = T(x) + T(z);$$

$$T(cx) = \lim_n T_n(cx) = \lim_n cT_n(x) = c \lim_n T_n(x) = cT(x).$$

Next, since (T_n) is a Cauchy sequence, $\{T_n : n \in \mathbb{N}^+\}$ is a bounded subset of $B(V, W)$. So there exists a finite constant $M \in \mathbb{R}^+$ with $\|T_n\| \leq M$ for all $n \in \mathbb{N}^+$. For any $x \in V$, $\|T(x)\|_W = \|\lim_n T_n(x)\|_W = \lim_n \|T_n(x)\|_W$, where $\|T_n(x)\|_W \leq \|T_n\| \cdot \|x\|_V \leq M \|x\|_V$ for all n . So $\|T(x)\|_W \leq M \|x\|_V$ for all $x \in V$, proving that T is bounded. Finally, we show $T_n \rightarrow T$ in $B(V, W)$. Given $\epsilon > 0$, we find n_0 so that $n \geq n_0$ implies $\|T_n - T\| < \epsilon$. Since (T_n) is a Cauchy sequence, we can choose n_0 so that $m, n \geq n_0$ implies $\|T_n - T_m\| < \epsilon/4$. Next, given $x \in V$ with $\|x\|_V = 1$, choose $m \geq n_0$ (depending on x) so that $\|T_m(x) - T(x)\|_W < \epsilon/4$. Now, for $n \geq n_0$,

$$\begin{aligned} \|(T_n - T)(x)\|_W &= \|T_n(x) - T_m(x) + T_m(x) - T(x)\|_W \\ &\leq \|(T_n - T_m)(x)\|_W + \|T_m(x) - T(x)\|_W < \epsilon/4 + \epsilon/4 = \epsilon/2. \end{aligned}$$

Note that n_0 is independent of x here, so we see that $\|T_n - T\| \leq \epsilon/2 < \epsilon$, completing the proof.

14.18 Dual Space of a Hilbert Space

In Chapter 13, we studied the *dual space* V^* of an F -vector space V , which is the vector space of all linear maps from V to F . For finite-dimensional V , we saw that V and V^* were isomorphic. In the case of a finite-dimensional real inner product space V , we could realize this isomorphism by mapping $y \in V$ to the linear functional $R_y \in V^*$ given by $R_y(x) = \langle x, y \rangle$ for $x \in V$. Similar results held for complex inner product spaces, but there we had to distinguish between linear maps and semi-linear maps.

Now let H be a Hilbert space. We define the *dual space* H^* to be the set $B(H, \mathbb{C})$ of all *continuous* linear maps from H to \mathbb{C} . A function $f : H \rightarrow \mathbb{C}$ is in H^* iff for all $x, y, z_n, z \in H$ and all $c \in \mathbb{C}$, $f(x+y) = f(x) + f(y)$, $f(cx) = cf(x)$, and $z_n \rightarrow z$ in H implies $f(z_n) \rightarrow f(z)$ in \mathbb{C} . For example, given $y \in H$, consider the map $R_y : H \rightarrow \mathbb{C}$ defined by $R_y(x) = \langle x, y \rangle$ for all $x \in H$. As noted in §14.13, each R_y is linear and continuous, so $R_y \in H^*$ for all $y \in H$.

Since \mathbb{C} is complete, it follows from the theorem proved in §14.17 that $H^* = B(H, \mathbb{C})$ is a Banach space with norm $\|f\| = \sup\{|f(x)| : x \in H, \|x\| = 1\}$ for $f \in H^*$. Our goal here is to define a semi-linear bijective isometry $R : H \rightarrow H^*$, which shows that the Hilbert space H and the normed vector space H^* are essentially isomorphic (up to conjugation of

scalars). We can use this semi-isomorphism to define an inner product on H^* that makes H^* a Hilbert space.

The map $R : H \rightarrow H^*$ is given by $R(y) = R_y$ for all $y \in H$, where $R_y(x) = \langle x, y \rangle$ for all $x \in H$. One checks, using the inner product axioms, that R is a semi-linear map. Next we prove that R is a bijection; in other words, *for every $f \in H^*$, there exists a unique $y \in H$ with $f = R_y$.* Fix $f \in H^*$. To prove uniqueness of y , assume $f = R_w = R_y$ for some $w, y \in H$; we show $w = y$. Compute

$$\|w - y\|^2 = \langle w - y, w - y \rangle = \langle w - y, w \rangle - \langle w - y, y \rangle = R_w(w - y) - R_y(w - y) = 0,$$

so $w - y = 0$ and $w = y$. To prove existence of y , consider the null space W of f , namely $W = \{x \in H : f(x) = 0\} = f^{-1}[\{0\}]$. W is a subspace of H , because it is the kernel of the linear map f . W is closed, because it is the inverse image of the closed set $\{0\}$ under the continuous map f . So we can write $H = W \oplus W^\perp$. If $W = H$, then f is the zero map, and we may take $y = 0$. If $W \neq H$, then there exists a nonzero $z \in W^\perp$. One readily checks that W^\perp is the one-dimensional space spanned by z (e.g., this follows from the fundamental homomorphism theorem for vector spaces). Let $y = cz$, where $c = \|z\|^{-2}f(z)$. We claim $f = R_y$. To prove this, fix $x \in H$, and write $x = w + dz$ for some $w \in W$ and $d \in \mathbb{C}$. On one hand, $f(x) = f(w + dz) = f(w) + df(z) = df(z)$. On the other hand, $R_y(x) = \langle w + dz, cz \rangle = \langle w, cz \rangle + \langle dz, cz \rangle = 0 + d\bar{c}\|z\|^2 = df(z)$. So the claim holds, and R is a bijection.

Finally, we check that R is an isometry. Fix $x, y \in H$ with $\|x\| = 1$ and $y \neq 0$. On one hand, by the Cauchy–Schwarz inequality,

$$|R_y(x)| = |\langle x, y \rangle| \leq \|y\| \cdot \|x\| = \|y\|,$$

so that $\|R_y\| \leq \|y\|$. On the other hand, letting $u = \|y\|^{-1}y$, we have $\|u\| = 1$ and

$$|R_y(u)| = |\langle y, u \rangle| / \|y\| = \|y\|,$$

so that $\|R_y\| \geq \|y\|$. Thus, $\|R_y\| = \|y\|$, and this equality also holds for $y = 0$. So R is indeed an isometry. It is now routine to check that H^* is a Hilbert space with inner product $\langle R_x, R_y \rangle = \langle y, x \rangle$ for $x, y \in H$; we reverse the order of x and y because R is semi-linear.

14.19 Adjoint

Next we discuss adjoints of operators on Hilbert spaces. One can approach this topic through the dual space H^* (as we did in §13.10), or as follows. An *operator* on a Hilbert space H is a continuous linear map $T : H \rightarrow H$; we write $B(H) = B(H, H)$ for the set of all such operators. We have seen that $B(H)$ is a Banach space (complete normed vector space). Given any $T \in B(H)$, we claim there exists a unique operator T^* on H , called the *adjoint* of T , such that $\langle T(x), y \rangle = \langle x, T^*(y) \rangle$ for all $x, y \in H$. Fix $y \in H$. The map f sending $x \in H$ to $\langle T(x), y \rangle$ is \mathbb{C} -linear, as one immediately verifies. It is also continuous, since $x_n \rightarrow x$ implies $T(x_n) \rightarrow T(x)$, hence $f(x_n) = \langle T(x_n), y \rangle \rightarrow \langle T(x), y \rangle = f(x)$. As we saw in §14.18, there exists a unique $w \in H$ with $f(x) = R_w(x) = \langle x, w \rangle$ for all $x \in H$. Writing $T^*(y) = w$ for each y , we obtain a unique function $T^* : H \rightarrow H$ satisfying $\langle T(x), y \rangle = \langle x, T^*(y) \rangle$ for all $x, y \in H$.

To continue, we need this lemma: if $u, v \in H$ satisfy $\langle x, u \rangle = \langle x, v \rangle$ for all $x \in H$, then $u = v$. For, taking $x = u - v$, we find that $\|u - v\|^2 = \langle u - v, u - v \rangle = \langle u - v, u \rangle - \langle u - v, v \rangle = 0$, forcing $u - v = 0$ and $u = v$.

Now we prove that T^* is a *linear* map. Given $y, z \in H$, is $T^*(y + z) = T^*(y) + T^*(z)$? Using the lemma, we answer this question in the affirmative by fixing $x \in H$ and computing

$$\begin{aligned}\langle x, T^*(y + z) \rangle &= \langle T(x), y + z \rangle = \langle T(x), y \rangle + \langle T(x), z \rangle = \langle x, T^*(y) \rangle + \langle x, T^*(z) \rangle \\ &= \langle x, T^*(y) + T^*(z) \rangle.\end{aligned}$$

Similarly, for $y \in H$ and $c \in \mathbb{C}$, $T^*(cy) = cT^*(y)$ holds because for each $x \in H$,

$$\langle x, T^*(cy) \rangle = \langle T(x), cy \rangle = \bar{c} \langle T(x), y \rangle = \bar{c} \langle x, T^*(y) \rangle = \langle x, cT^*(y) \rangle.$$

Finally, is T^* *continuous*? Fix $y \in H$, and use the Cauchy–Schwarz inequality and the boundedness of T to compute

$$\|T^*(y)\|^2 = \langle T^*(y), T^*(y) \rangle = \langle T(T^*(y)), y \rangle \leq \|T(T^*(y))\| \cdot \|y\| \leq \|T\| \cdot \|T^*(y)\| \cdot \|y\|.$$

If $\|T^*(y)\| > 0$, we divide by $\|T^*(y)\|$ to see that $\|T^*(y)\| \leq \|T\| \cdot \|y\|$; and this inequality also holds if $\|T^*(y)\| = 0$. It now follows that T^* is bounded, hence continuous. In fact, the proof shows that $\|T^*\| \leq \|T\|$.

Here are some properties of adjoint operators: for $S, T \in B(H)$ and $c \in \mathbb{C}$, $(S + T)^* = S^* + T^*$, $(cS)^* = \bar{c}(S^*)$, $S^{**} = S$, $(S \circ T)^* = T^* \circ S^*$, and $\|T^*\| = \|T\|$. We can prove these using the lemma stated above. For instance, $S^{**} = S$ since for all $x, y \in H$,

$$\langle x, S^{**}(y) \rangle = \langle S^*(x), y \rangle = \overline{\langle y, S^*(x) \rangle} = \overline{\langle S(y), x \rangle} = \langle x, S(y) \rangle.$$

Similarly, $(S \circ T)^* = (T^* \circ S^*)$ because for all $x, y \in H$,

$$\begin{aligned}\langle x, (S \circ T)^*(y) \rangle &= \langle (S \circ T)(x), y \rangle = \langle S(T(x)), y \rangle = \langle T(x), S^*(y) \rangle = \langle x, T^*(S^*(y)) \rangle \\ &= \langle x, (T^* \circ S^*)(y) \rangle.\end{aligned}$$

We leave semi-linearity of the map $S \mapsto S^*$ as an exercise. Finally, we showed earlier that $\|T^*\| \leq \|T\|$ for all $T \in B(H)$. Replacing T by T^* gives $\|T\| = \|T^{**}\| \leq \|T^*\|$, so $\|T^*\| = \|T\|$.

Now that we have the concept of an adjoint operator, we can generalize the special types of matrices studied in Chapter 7 to the setting of Hilbert spaces. An operator T on a Hilbert space H is called *self-adjoint* iff $T^* = T$; *positive* iff $T^* = T$ and $\langle T(x), x \rangle \geq 0$ for all $x \in H$; *normal* iff $T \circ T^* = T^* \circ T$; and *unitary* iff $T \circ T^* = \text{id}_H = T^* \circ T$. We indicate some basic properties of these operators in the exercises, but we must refer the reader to more advanced texts for a detailed account of the structure theory of normal operators on Hilbert spaces.

14.20 Summary

- Definitions for Metric Spaces.* Table 14.1 summarizes definitions of concepts related to metric spaces, sequences, and continuous functions.
- Examples of Metric Spaces.* \mathbb{R}^m and \mathbb{C}^m are metric spaces with the Euclidean metric $d_2(x, y) = \sqrt{\sum_{k=1}^m |x_k - y_k|^2}$. Any set X has the discrete metric $d(x, y) = 0$ for $x = y$, $d(x, y) = 1$ for $x \neq y$. A product $X = X_1 \times \cdots \times X_m$ of metric spaces X_k has metrics $d_1(x, y) = \sum_{k=1}^m d_{X_k}(x_k, y_k)$ and $d_\infty(x, y) = \max\{d_{X_k}(x_k, y_k) : 1 \leq k \leq m\}$.

TABLE 14.1

Definitions of Concepts for Metric Spaces.

Concept	Definition
metric space (X, d)	$\forall x, y \in X, 0 \leq d(x, y) < \infty; d(x, y) = 0 \Leftrightarrow x = y$ $\forall x, y \in X, d(x, y) = d(y, x)$ (symmetry) $\forall x, y, z \in X, d(x, z) \leq d(x, y) + d(y, z)$ (triangle ineq.)
convergent sequence (x_n)	$x_n \rightarrow x \Leftrightarrow \forall \epsilon > 0, \exists n_0 \in \mathbb{N}, \forall n \geq n_0, d(x_n, x) < \epsilon$.
Cauchy sequence (x_n)	$\forall \epsilon > 0, \exists n_0 \in \mathbb{N}, \forall m, n \geq n_0, d(x_n, x_m) < \epsilon$.
closed set C	If $x_n \rightarrow x$ and all $x_n \in C$, then $x \in C$.
open set U	$\forall x \in U, \exists r > 0, \forall y \in X, d(x, y) < r \Rightarrow y \in U$.
bounded set S	$\exists x \in X, \exists M \in \mathbb{R}, \forall y \in S, d(x, y) \leq M$.
totally bounded set S (Exc. 45)	$\forall \epsilon > 0, \exists m \in \mathbb{N}^+, \exists x_1, \dots, x_m \in X, S \subseteq \bigcup_{k=1}^m B(x_k; \epsilon)$.
sequentially compact set K	Any sequence (x_n) in K has a subsequence converging to a point of K .
topologically compact set K	Whenever $K \subseteq \bigcup_{i \in I} U_i$ with all U_i open, there is a finite $F \subseteq I$ with $K \subseteq \bigcup_{i \in F} U_i$.
complete set K	Every Cauchy sequence in K converges to a point of K .
continuous $f : X \rightarrow Y$	Whenever $x_n \rightarrow x$ in X , $f(x_n) \rightarrow f(x)$ in Y .

3. *Convergent Sequences and Cauchy Sequences.* The limit of a convergent sequence is unique. Convergent sequences are Cauchy sequences; the converse holds in complete spaces. Convergent sequences and Cauchy sequences are bounded. Every subsequence of a convergent sequence converges to the same limit as the full sequence. If one subsequence of a Cauchy sequence converges, then the full sequence converges to the same limit.
4. *Closed Sets.* In any metric space X , \emptyset and X are closed. Finite subsets of X are closed. The union of finitely many closed sets is closed. The intersection of arbitrarily many closed sets is closed. C is closed iff $X \sim C$ is open. Closed intervals $[a, b]$ in \mathbb{R} are closed sets. If C_1, \dots, C_m are closed in X_1, \dots, X_m respectively, then $C_1 \times \dots \times C_m$ is closed in the product metric space $X_1 \times \dots \times X_m$ (with metric d_1 or d_∞).
5. *Open Sets.* In any metric space X , \emptyset and X are open. The union of arbitrarily many open sets is open. The intersection of finitely many open sets is open. U is open iff $X \sim U$ is closed. Open balls $B(x; r)$ are open sets. Open intervals (a, b) in \mathbb{R} are open sets, and every open set in \mathbb{R} is a disjoint union of countably many open intervals. If U_1, \dots, U_m are open in X_1, \dots, X_m respectively, then $U_1 \times \dots \times U_m$ is open in the product metric space $X_1 \times \dots \times X_m$ (with metric d_1 or d_∞).
6. *Continuous Functions.* A map $f : X \rightarrow Y$ is continuous iff $x_n \rightarrow x$ in X implies $f(x_n) \rightarrow f(x)$ in Y iff for all open $V \subseteq Y$, $f^{-1}[V]$ is open in X iff for all closed $D \subseteq Y$, $f^{-1}[D]$ is closed in X iff for all $\epsilon > 0$ and $x \in X$, there exists $\delta > 0$ such that for all $z \in X$ with $d_X(x, z) < \delta$, $d_Y(f(x), f(z)) < \epsilon$. If K is compact in X and f is continuous, then $f[K]$ is compact in Y . The direct image of an open (or closed) subset of X under a continuous map need not be open (or closed) in Y .
7. *Compact Sets.* A subset K of a metric space X is sequentially compact iff every sequence of points in K has a subsequence converging to a point of K . In metric

spaces, sequential compactness is equivalent to topological compactness (every open cover of K has a finite subcover). Finite subsets of X are compact. Finite unions and arbitrary intersections of compact subsets are compact. Compact subsets of X must be closed in X and bounded, but the converse does not hold in general. However, in \mathbb{R}^m and \mathbb{C}^m with the metrics d_1 , d_2 , and d_∞ , a subset K is compact iff K is closed and bounded. A closed subset of a compact set is compact. The direct image of a compact set under a continuous function is compact; but the inverse image of a compact set may not be compact. K is compact iff K is complete and totally bounded.

8. *Complete Spaces.* A metric space X is complete iff every Cauchy sequence in X converges to a point of X . Compactness implies completeness, but not conversely. \mathbb{R}^m and \mathbb{C}^m with the metrics d_1 , d_2 , and d_∞ are complete but not compact. A subset of a complete space is complete iff it is closed.
9. *Hilbert Spaces.* A Hilbert space is a complex inner product space, with norm $\|x\| = \sqrt{\langle x, x \rangle}$ and metric $d(x, y) = \|x - y\|$, that is complete as a metric space. \mathbb{C}^m is an m -dimensional Hilbert space. For any set X , the set $\ell_2(X)$ of functions $f : X \rightarrow \mathbb{C}$ with $\sum_{x \in X} |f(x)|^2 < \infty$ is a Hilbert space with inner product $\langle f, g \rangle = \sum_{x \in X} f(x)\overline{g(x)}$ and norm $\|f\| = \sqrt{\sum_{x \in X} |f(x)|^2}$. The set of integrable functions on a measure space (X, μ) is a Hilbert space with inner product $\langle f, g \rangle = \int_X f(x)\overline{g(x)} d\mu$ and norm $\|f\| = \sqrt{\int_X |f(x)|^2 d\mu}$.
10. *Properties of Hilbert Spaces.* All x, y in a Hilbert space H satisfy the *Cauchy–Schwarz inequality*: $|\langle x, y \rangle| \leq \|x\| \cdot \|y\|$; the *parallelogram law*:

$$\|x + y\|^2 + \|x - y\|^2 = 2\|x\|^2 + 2\|y\|^2;$$

and the *polarization identity*:

$$\|x + y\|^2 - \|x - y\|^2 + i\|x + iy\|^2 - i\|x - iy\|^2 = 4\langle x, y \rangle.$$

If x and y are orthogonal (meaning $\langle x, y \rangle = 0$), then the *Pythagorean theorem* holds: $\|x \pm y\|^2 = \|x\|^2 + \|y\|^2$. For pairwise orthogonal $x_1, \dots, x_m \in H$ and $c_1, \dots, c_m \in \mathbb{C}$, $\|c_1x_1 + \dots + c_mx_m\|^2 = |c_1|^2\|x_1\|^2 + \dots + |c_m|^2\|x_m\|^2$. The Hilbert space operations are continuous: if $x_n \rightarrow x$ and $y_n \rightarrow y$ in H and $c_n \rightarrow c$ in \mathbb{C} , then $x_n + y_n \rightarrow x + y$, $c_nx_n \rightarrow cx$, $\langle x_n, y_n \rangle \rightarrow \langle x, y \rangle$, and $\|x_n\| \rightarrow \|x\|$.

11. *Closed Convex Sets.* Let C be a nonempty, closed, convex subset of a Hilbert space. C has a unique element of minimum norm. For each $w \in H$, there exists a unique point of C closest to w .
12. *Orthogonal Complements.* For each subset S of a Hilbert space H , the orthogonal complement $S^\perp = \{v \in H : \langle v, w \rangle = 0 \text{ for all } w \in S\}$ is a closed subspace of H . The map $W \mapsto W^\perp$ is an order-reversing bijection on the lattice of all closed subspaces of H . For each closed subspace W , $H = W \oplus W^\perp$, $W^{\perp\perp} = W$, and in the unique expression of $x \in H$ as a sum of $y \in W$ and $z \in W^\perp$, y (resp. z) is the closest point to x in W (resp. W^\perp). For any subset S , $S^{\perp\perp}$ is the smallest closed subspace containing S .
13. *Orthonormal Sets.* A subset X of a Hilbert space H is orthonormal iff $\langle x, x \rangle = 1$ and $\langle x, y \rangle = 0$ for all $x \neq y$ in X . For X orthonormal and $w \in H$, Bessel's inequality states that $\sum_{x \in X} |\langle w, x \rangle|^2 \leq \|w\|^2$. Equality holds for all $w \in H$ iff X is a maximal orthonormal set.

14. *The Isomorphism $H \cong \ell_2(X)$.* Zorn's Lemma assures that every Hilbert space H has a maximal orthonormal subset X . The map $f : H \rightarrow \ell_2(X)$ sending $w \in H$ to the function $f_w : X \rightarrow \mathbb{C}$, given by $f_w(x) = \langle w, x \rangle$ for $x \in X$, is a bijective, continuous, linear isometry (Hilbert space isomorphism). It follows that $\sum_{x \in X} |\langle w, x \rangle|^2 = \|w\|^2$ and $\sum_{x \in X} \langle w, x \rangle \overline{\langle z, x \rangle} = \langle w, z \rangle$ for $w, z \in H$ (Parseval's identity).
15. *Continuous Linear Maps.* A linear map $T : V \rightarrow W$ between two normed vector spaces is continuous ($v_n \rightarrow v$ in V implies $T(v_n) \rightarrow T(v)$ in W) iff T is continuous at zero ($v_n \rightarrow 0$ in V implies $T(v_n) \rightarrow 0$ in W) iff T is bounded (for some finite M , $\|T(v)\| \leq M\|v\|$ for all $v \in V$). The set $B(V, W)$ of continuous linear maps from V to W is a normed vector space with norm $\|T\| = \sup\{\|T(x)\| : x \in V, \|x\| = 1\}$. If W is complete, then $B(V, W)$ is complete.
16. *Dual of a Hilbert Space.* Given a Hilbert space H , the dual space H^* is $B(H, \mathbb{C})$, the set of bounded (continuous) linear maps from H to \mathbb{C} . For every $f \in H^*$, there exists a unique $y \in H$ such that $f = R_y$, where $R_y(x) = \langle x, y \rangle$ for $x \in H$. The map $R : H \rightarrow H^*$ is a bijective semi-linear isometry. So, every Hilbert space is semi-isomorphic to its dual space.
17. *Adjoint Operators.* For each operator $T \in B(H) = B(H, H)$, there exists a unique adjoint operator $T^* \in B(H)$ satisfying $\langle T(x), y \rangle = \langle x, T^*(y) \rangle$ for all $x, y \in H$. For $S, T \in B(H)$ and $c \in \mathbb{C}$, $(S + T)^* = S^* + T^*$, $(cS)^* = \bar{c}(S^*)$, $S^{**} = S$, $(S \circ T)^* = T^* \circ S^*$, and $\|T^*\| = \|T\|$. The operator T is *self-adjoint* iff $T = T^*$; *positive* iff $T = T^*$ and $\langle T(x), x \rangle \geq 0$ for all $x \in H$; *normal* iff $T \circ T^* = T^* \circ T$; and *unitary* iff $T \circ T^* = \text{id}_H = T^* \circ T$.

14.21 Exercises

Unless otherwise specified, assume X and Y are arbitrary metric spaces and H is an arbitrary Hilbert space in these exercises.

- Which of the following functions define metrics on \mathbb{R}^2 ? Explain.
 - $d((x, y), (u, v)) = |x - u|$. (b) $d((x, y), (u, v)) = (x - u)^2 + (y - v)^2$.
 - $d((x, y), (u, v)) = \sqrt[3]{(x - u)^3 + (y - v)^3}$. (d) $d((x, y), (u, v)) = |x| + |y|$ if $(x, y) \neq (u, v)$, and 0 otherwise. (e) $d((x, y), (u, v)) = 3$ if $(x, y) \neq (u, v)$, and 0 otherwise.
- (a) Let X be the set of all bounded functions $f : [0, 1] \rightarrow \mathbb{R}$. Show that $d(f, g) = \sup\{|f(x) - g(x)| : x \in [0, 1]\}$ for $f, g \in X$ defines a metric on X . (b) Let Y be the set of all continuous functions $f : [0, 1] \rightarrow \mathbb{R}$. Show that $d(f, g) = \int_0^1 |f(x) - g(x)| dx$ for $f, g \in Y$ defines a metric on Y . (c) What is $d(x^2, x^3)$ using the metric in (a)? What is $d(x^2, x^3)$ using the metric in (b)? (d) Let Z be the set of all Riemann-integrable functions $f : [0, 1] \rightarrow \mathbb{R}$. Is d (defined as in (b)) a metric on Z ?
- Let p be a fixed prime integer. Define $d_p : \mathbb{Q} \times \mathbb{Q} \rightarrow \mathbb{R}$ by letting $d_p(x, y) = 0$ if $x = y$, $d_p(x, y) = 1/p^k$ if $x - y$ is a nonzero integer and p^k is the largest power of p dividing $x - y$, and $d_p(x, y) = 1$ otherwise. (a) Prove: for all $x, y, z \in \mathbb{Q}$, $d_p(x, z) \leq \max(d_p(x, y), d_p(y, z))$. (b) Prove that (\mathbb{Q}, d_p) is a metric space.
- Equivalent Metrics.* We say that two metrics d and d' defined on the same set

X are *equivalent* iff they have the same open sets, i.e., for all $U \subseteq X$, U is open in the metric space (X, d) iff U is open in the metric space (X, d') . (a) Prove that equivalent metrics have the same closed sets, the same compact sets, and the same convergent sequences. (b) Prove that $f : X \rightarrow Y$ is continuous relative to d iff f is continuous relative to d' . (c) Show that equivalence of metrics is an equivalence relation on the set of all metrics on a fixed set X .

5. Let (X, d) be a metric space. Fix $m > 0$, and define $d' : X \times X \rightarrow \mathbb{R}$ by $d'(x, y) = \min(m, d(x, y))$ for $x, y \in X$. (a) Prove that d' is a metric on X . (b) Prove that d' is equivalent to d (see Exercise 4). (c) Show that all subsets of X are bounded relative to d' . (d) Use (b) and (c) to construct an example of a closed, bounded, non-compact subset of a metric space.
6. *Comparable Metrics.* We say that two metrics d and d' defined on the same set X are *comparable* iff there exist positive, finite constants M, N such that for all $x, y \in X$, $d(x, y) \leq M d'(x, y)$ and $d'(x, y) \leq N d(x, y)$. (a) Show that if d and d' are comparable, then they are equivalent (see Exercise 4). (b) Show that d_2 and d_∞ are comparable metrics on \mathbb{R}^n and \mathbb{C}^n . (c) Show that d_1 and d_∞ are comparable metrics on \mathbb{R}^n and \mathbb{C}^n . (d) Show that comparability is an equivalence relation on the set of all metrics on a fixed set X . (e) Give an example of two equivalent metrics on \mathbb{R} that are not comparable.
7. Let X have metric d , let Z be any set, and let $f : X \rightarrow Z$ be a bijection. Show there exists a unique metric on Z such that f is an isometry.
8. (a) Give an example of two equivalent metrics on \mathbb{R} (see Exercise 4) such that \mathbb{R} is complete with respect to one of the metrics but not the other. (Study $f : (-\pi/2, \pi/2) \rightarrow \mathbb{R}$ given by $f(x) = \tan x$.) (b) Given two comparable metrics d_1 and d_2 on a set X (see Exercise 6), show that (X, d_1) is complete iff (X, d_2) is complete.
9. (a) Verify the metric space axioms for d' , d_1 , and d_∞ (defined in §14.1). (b) Given metric spaces X_1, \dots, X_m , show that the metrics d_1 and d_∞ on $X = X_1 \times \dots \times X_m$ are comparable.
10. (a) Negate the definition of convergent sequence to obtain the formal definition of a non-convergent sequence. (b) Use (a) to prove carefully that the sequence $(n : n \in \mathbb{N})$ is non-convergent in the metric space \mathbb{R} .
11. A sequence $(x_n : n \in \mathbb{N})$ is called *eventually constant* iff there exists $n_0 \in \mathbb{N}$ such that for all $n \geq n_0$, $x_n = x_{n_0}$. (a) Prove every eventually constant sequence in (X, d) converges to x_{n_0} . (b) If d is the discrete metric on X , prove that every convergent sequence is eventually constant.
12. Fix $k \in \mathbb{N}$. (a) Show that a sequence $(x_n : n \in \mathbb{N})$ converges to x iff the “tail” $(x_n : n \geq k) = (x_{n+k} : n \in \mathbb{N})$ converges to x . (b) Show that $(x_n : n \in \mathbb{N})$ is Cauchy iff $(x_n : n \geq k)$ is Cauchy. (c) Show that $(x_n : n \in \mathbb{N})$ is bounded iff $(x_n : n \geq k)$ is bounded.
13. (a) Show that if $(x_n : n \in \mathbb{N})$ is a convergent sequence in X , then $\{x_n : n \in \mathbb{N}\}$ is a bounded subset of X . (b) Show that if $(x_n : n \in \mathbb{N})$ is a Cauchy sequence in X , then $\{x_n : n \in \mathbb{N}\}$ is a bounded subset of X .
14. Let $f : X \rightarrow Y$ be a continuous map between metric spaces, and let (x_n) be a sequence in X . (a) If (x_n) is convergent, or Cauchy, or bounded, must the same be true of the sequence $(f(x_n))$? Explain. (b) If $(f(x_n))$ is convergent, or Cauchy, or bounded, must the same be true of (x_n) ? Explain.

15. Let $(x_n) = ((x_n(1), x_n(2), \dots, x_n(m)) : n \in \mathbb{N})$ be a sequence in a product metric space $X = X_1 \times X_2 \times \dots \times X_m$. (a) Prove (x_n) converges to $y = (y_1, \dots, y_m)$ in (X, d_1) iff $(x_n(k))$ converges to y_k for $1 \leq k \leq m$. (b) Prove (x_n) is Cauchy iff $(x_n(k))$ is Cauchy for $1 \leq k \leq m$. (c) Repeat (a) and (b) for (X, d_∞) . (d) Repeat (a) and (b) for $X = \mathbb{R}^m$ or \mathbb{C}^m using the Euclidean metric d_2 .
16. Let $(x_n : n \in \mathbb{N})$ be a sequence in (X, d) . Show that the set S of all limits of subsequences of $(x_n : n \in \mathbb{N})$ is a closed subset of X .
17. (a) Give an example of a sequence $(x_n : n \in \mathbb{N})$ in \mathbb{R} such that every $z \in \mathbb{Z}$ is a limit of some subsequence of (x_n) . (b) Give an example of a sequence $(x_n : n \in \mathbb{N})$ in \mathbb{R} such that every $r \in [0, 1]$ is a limit of some subsequence of (x_n) .
18. Show that every subset of a discrete metric space is both open and closed.
19. Decide (with proof) whether each subset of (\mathbb{R}^2, d_2) is closed.
 (a) $\{(x, y) \in \mathbb{R}^2 : y \geq 0\}$; (b) $\{(x, y) \in \mathbb{R}^2 : x^2 + y^2 < 1\}$; (c) $\mathbb{Z} \times \mathbb{Z}$; (d) the graph $\{(x, f(x)) : x \in \mathbb{R}\}$ where $f : \mathbb{R} \rightarrow \mathbb{R}$ is continuous.
20. Let $Z \subseteq Y \subseteq X$. (a) Show that if Z is closed in the subspace Y , and Y is closed in X , then Z is closed in X . (b) Give an example to show that (a) can fail if Y is not closed in X . (c) Show that if Z is open in the subspace Y , and Y is open in X , then Z is open in X . (d) Give an example to show that (c) can fail if Y is not open in X .
21. *Closure of a Set.* Given a subset S of the metric space (X, d) , the *closure of S in X* , denoted \overline{S} , is the intersection of all closed subsets of X containing S .
 (a) Show that $S \subseteq \overline{S}$, \overline{S} is a closed set, and for any closed set T , if $S \subseteq T$ then $\overline{S} \subseteq T$. (b) Show that \overline{S} is the set S' of all $y \in X$ such that there exists a sequence (x_n) with all $x_n \in S$ and $\lim_{n \rightarrow \infty} x_n = y$. (First show S' is closed.)
 (c) In (\mathbb{R}, d_2) , compute the closure of these sets: $S = (a, b)$, where $a < b$; $S = \mathbb{Q}$; $S = \{1/n : n \in \mathbb{N}^+\}$; $S = \mathbb{Z}$.
22. *Boundary of a Set.* Given a subset S of the metric space (X, d) , the *boundary of S* is the set $\text{Bdy}(S)$ consisting of all $x \in X$ such that for all $\epsilon > 0$, $B(x; \epsilon)$ contains at least one point in S and at least one point not in S . (a) Prove $\text{Bdy}(S) = \overline{S} \cap \overline{X \sim S}$ (see Exercise 21). (b) Prove $\text{Bdy}(S)$ is a closed set. (c) In (\mathbb{R}, d_2) , compute the boundary of these sets: $S = (a, b)$, where $a < b$; $S = \mathbb{Q}$; $S = \mathbb{R}$; $S = \mathbb{Z}$. (d) In (\mathbb{R}^2, d_2) , compute the boundary of these sets: $S = \mathbb{R} \times \{0\}$; $S = \{(x, y) : x^2 + y^2 < 1\}$; $S = [0, 1] \times \mathbb{Q}$.
23. Let $X = X_1 \times \dots \times X_m$ be a product of metric spaces. (a) Prove: If S_i is closed in X_i for all i , then $S_1 \times \dots \times S_m$ is closed in (X, d_1) . (b) Prove: If S_i is open in X_i for all i , then $S_1 \times \dots \times S_m$ is open in (X, d_1) . (c) Prove: If S_i is compact in X_i for all i , then $S_1 \times \dots \times S_m$ is compact in (X, d_1) . (d) Do the results in (a), (b), and (c) hold for (X, d_∞) ? Do they hold when $X = (\mathbb{R}^m, d_2)$ and each $X_i = (\mathbb{R}, d_2)$?
24. Let W be a subspace of a normed vector space V . Show that the closure \overline{W} (see Exercise 21) is a closed subspace of V .
25. Sketch the open ball $B((2, 1); 1)$ in \mathbb{R}^2 for each of the metrics d_1 , d_2 , and d_∞ .
26. *Interior of a Set.* Given a subset S of the metric space (X, d) , the *interior of S in X* , denoted $\text{Int}(S)$ or S° , is the union of all open subsets of X contained in S . (a) Show that $\text{Int}(S) \subseteq S$, $\text{Int}(S)$ is an open set, and for any open set U , if $U \subseteq S$ then $U \subseteq \text{Int}(S)$. (b) Show that $\text{Int}(S) = X \sim \overline{X \sim S}$. (c) In (\mathbb{R}, d_2) , give an example of a set S with $\text{Int}(S) = \emptyset$ and $\overline{S} = \mathbb{R}$.

27. Fix $x \in X$ and $r > 0$. (a) Show that the *closed ball* $B[x; r] = \{y \in X : d(x, y) \leq r\}$ is a closed subset of X . (b) Show that open and closed balls in normed vector spaces are convex. (c) Give an example to show that $\overline{B(x; r)}$ can be a proper subset of $B[x; r]$.
28. (a) Given $S \subseteq X$ and $r > 0$, let $N_r(S) = \{y \in X : d(x, y) < r \text{ for some } x \in S\}$. Show that $N_r(S)$ is always an open set.
 (b) Let $N_r[S] = \{y \in X : d(x, y) \leq r \text{ for some } x \in S\}$. Must $N_r[S]$ be a closed set? Explain.
29. Show that every open set in \mathbb{R} is the union of countably many disjoint open intervals.
30. (a) Prove that every isometry is continuous. (b) A *contraction* is a function $f : X \rightarrow Y$ such that $d_Y(f(x_1), f(x_2)) \leq d_X(x_1, x_2)$ for all $x_1, x_2 \in X$. Prove every contraction is continuous. (c) Deduce from (a) and (b) that id_X is continuous, and constant functions are continuous.
31. (a) Let $f : X \rightarrow Y$ be an isometry. Prove: for all $x_n \in X$, (x_n) is a Cauchy sequence in X iff $(f(x_n))$ is a Cauchy sequence in Y . (b) Does (a) hold if f is a contraction?
32. Let $X = X_1 \times \cdots \times X_m$ be a product metric space using the metric d_1 defined in §14.1. (a) Show that $p_i : X \rightarrow X_i$, given by $p_i((x_1, \dots, x_m)) = x_i$ for $(x_1, \dots, x_m) \in X$, is continuous for $1 \leq i \leq m$. (b) Show that $g : Y \rightarrow X$ is continuous iff $p_i \circ g$ is continuous for $1 \leq i \leq m$. (c) Explain why (a) and (b) also hold for (X, d_∞) and (\mathbb{C}^m, d_2) .
33. Let $Z = \{z \in \mathbb{C} : |z| = 1\}$ with the Euclidean metric. Define $f : [0, 2\pi) \rightarrow Z$ by $f(x) = (\cos x, \sin x)$ for $x \in [0, 2\pi)$. Verify that f is a continuous bijection, but f^{-1} is not continuous.
34. (a) Prove that if X has the discrete metric, then every function $f : X \rightarrow Y$ is continuous. (b) Use (a) to give an example of a continuous bijection whose inverse is not continuous.
35. Let V, W, Z be normed vector spaces. Suppose (T_n) is a sequence in $B(V, W)$ converging to T , and (U_n) is a sequence in $B(W, Z)$ converging to U . Prove $U_n \circ T_n \rightarrow U \circ T$ in $B(V, Z)$.
36. Prove $f : X \rightarrow Y$ is continuous iff for all $x_0 \in X$ and all $\epsilon > 0$, there exists $\delta > 0$ such that for all $x \in X$, if $d_X(x, x_0) < \delta$, then $d_Y(f(x), f(x_0)) < \epsilon$.
37. *Uniform Continuity.* A function $f : X \rightarrow Y$ is called *uniformly continuous* iff for all $\epsilon > 0$, there exists $\delta > 0$ such that for all $x_1, x_2 \in X$, if $d_X(x_1, x_2) < \delta$, then $d_Y(f(x_1), f(x_2)) < \epsilon$. (Note that the δ appearing here depends on ϵ but not on x_1 and x_2 , whereas the δ in Exercise 36 depends on both ϵ and x_0). (a) Prove that contractions and isometries are uniformly continuous. (b) Prove that bounded linear maps on normed vector spaces are uniformly continuous. (c) Prove that addition (viewed as a map from $(H \times H, d_1)$ to H) is uniformly continuous. (d) Show that the norm on a normed vector space is uniformly continuous. (e) Show that $f : \mathbb{R} \rightarrow \mathbb{R}$, given by $f(x) = x^2$ for $x \in \mathbb{R}$, is continuous but not uniformly continuous.
38. (a) Prove that the union of a finite collection of compact subsets of X is compact.
 (b) Show that a subset S of a compact set is compact iff S is closed. (c) Prove that the intersection of any collection of compact subsets of X is compact.

39. Suppose X is a discrete metric space. (a) Under what conditions is X compact? (b) Under what conditions is X complete?
40. Let $f : X \rightarrow Y$ be continuous. Answer each question below with a proof or a counterexample. (a) If U is open in X , must $f[U]$ be open in Y ? (b) If C is closed in X , must $f[C]$ be closed in Y ? (c) If K is compact in Y , must $f^{-1}[K]$ be compact in X ?
41. *Theorem:* For all $a < b$ in \mathbb{R} , the closed interval $[a, b]$ is (sequentially) compact. Prove this theorem by completing the following outline. Let $\mathbf{x} = (x_n : n \in \mathbb{N})$ be a sequence of points in $[a, b]$. We will construct a sequence of nested closed intervals $I_0, I_1, \dots, I_k, \dots$, where $I_0 = [a, b]$, $I_{k+1} \subseteq I_k$ for all $k \geq 0$, and the length of I_k is $(b - a)/2^k$ for all $k \geq 0$. I_{k+1} will be either the left half or the right half of I_k . We will also construct a list of sequences $\mathbf{x}^0 = \mathbf{x}, \mathbf{x}^1, \mathbf{x}^2, \dots$, such that each sequence \mathbf{x}^{k+1} is a subsequence of the previous sequence \mathbf{x}^k , and for all k , all terms of \mathbf{x}^k lie in I_k . The construction proceeds recursively, with I_0 and \mathbf{x}^0 given above. Suppose that I_k and \mathbf{x}^k have been constructed with the stated properties. Each term in \mathbf{x}^k lies in either the left half or the right half of I_k . So one of these two halves, call it I_{k+1} , must contain infinitely many terms of the sequence \mathbf{x}^k . Define \mathbf{x}^{k+1} to be the subsequence of \mathbf{x}^k consisting of all terms of \mathbf{x}^k that lie in I_{k+1} . (a) Verify by induction that \mathbf{x}^k and I_k have the properties stated above. (b) Use the fact that every bounded subset of \mathbb{R} has a greatest lower bound and least upper bound to see that $\bigcap_{k \geq 0} I_k$ must consist of a single point $x^* \in \mathbb{R}$. (c) Consider the diagonal sequence (y_k) , where $y_k = \mathbf{x}_k^k$. Check that this sequence is a subsequence of the original sequence (x_n) , and $y_k \in I_k$ for all k . (d) Show that (y_k) converges to x^* .
42. Use Exercise 41 and Exercise 23 to prove that a subset K of \mathbb{R}^m (using any of the metrics d_1 , d_2 , or d_∞) is compact iff K is closed and bounded.
43. *Extreme Value Theorem.* Prove: if X is a compact metric space and $f : X \rightarrow \mathbb{R}$ is continuous, then there exist $x_1, x_2 \in X$ such that for all $x \in X$, $f(x_1) \leq f(x) \leq f(x_2)$. In other words, a continuous real-valued function on a compact domain attains its maximum and minimum value somewhere on that domain.
44. Let c_n be positive real constants with $\sum_{n=1}^{\infty} c_n^2 < \infty$. Show that $K = \{f \in \ell_2(\mathbb{N}^+) : 0 \leq |f(n)| \leq c_n \text{ for all } n\}$ is a compact subset of the Hilbert space $\ell_2(\mathbb{N}^+)$. (Recycle some ideas from the proof in Exercise 41.)
45. Call a subset S of X *totally bounded* iff for all $\epsilon > 0$, there exist finitely many points $x_1, \dots, x_m \in X$ such that $S \subseteq \bigcup_{i=1}^m B(x_i; \epsilon)$. Prove that S is (sequentially) compact iff S is complete and totally bounded.
46. Prove that a topologically compact metric space X is sequentially compact. (If X is not sequentially compact, fix a sequence (x_n) in X with no convergent subsequence. Without loss of generality, assume no two x_n 's are equal. For each $x \in X$, prove that there is an open ball $U_x = B(x; \epsilon(x))$ containing x such that at most one x_n lies in U_x . Show that $\{U_x : x \in X\}$ is an open cover of X with no finite subcover.)
47. A metric space X is called *separable* iff there exists a countable set $S = \{x_n : n \in \mathbb{N}^+\}$ with $\overline{S} = X$. (a) Prove that sequentially compact metric spaces are separable. (b) Prove that given any open cover $\{U_i : i \in I\}$ of a separable metric space X , there exists a countable subcover $\{U_{i_m} : m \in \mathbb{N}\}$ with $X = \bigcup_{m=0}^{\infty} U_{i_m}$. (Fix a set U_{i_0} in the open cover. Given S as above, for

each $n, k \in \mathbb{N}^+$ let $V_{n,k}$ be one particular set U_i in the open cover such that $B(x_n; 1/k) \subseteq U_i$, if such a set exists; otherwise let $V_{n,k} = U_{i_0}$. Show that the set of all $V_{n,k}$'s gives a countable subcover.)

48. Prove that a sequentially compact metric space X is topologically compact. (By Exercise 47, reduce to showing that every countably infinite open cover $\{V_n : n \in \mathbb{N}^+\}$ of X has a finite subcover. If not, then all of the closed sets $C_n = X \sim (V_1 \cup \dots \cup V_n)$ are nonempty; pick $x_n \in C_n$ and study the sequence (x_n) .)
49. Let X and Y be complete metric spaces with subsets $X_1 \subseteq X$ and $Y_1 \subseteq Y$. Assume $f_1 : X_1 \rightarrow Y_1$ is an isometry, $\overline{X_1} = X$, and $\overline{Y_1} = Y$. Show that there exists a unique extension of f_1 to an isometry $f : X \rightarrow Y$.
50. (a) Find necessary and sufficient conditions on $v, w \in H$ for equality to hold in the Cauchy–Schwarz inequality. (b) Find necessary and sufficient conditions on $v, w \in H$ for equality to hold in the triangle inequality for norms.
51. Show that every finite-dimensional complex inner product space is automatically complete relative to the metric induced from the inner product. (Use orthonormal bases to show that such a space is isometrically isomorphic to \mathbb{C}^n with the Euclidean metric.)
52. Let W be the subset of $H = \ell_2(\mathbb{N})$ consisting of all sequences with only finitely many nonzero terms. Show that W is a subspace of H that is not closed.
53. Let H be the Hilbert space of (Lebesgue-measurable) functions $f : [0, 1] \rightarrow \mathbb{C}$ such that $\int_0^1 |f(x)|^2 dx < \infty$. Show that the set of continuous $f \in H$ is a subspace of H that is not closed.
54. Let S be any subset of H . (a) Let $\langle S \rangle$ be the span of S , consisting of all finite \mathbb{C} -linear combinations of elements of S . Prove $S^\perp = \langle S \rangle^\perp$. (b) Let \overline{S} be the closure of S (see Exercise 21). Prove $S^\perp = (\overline{S})^\perp$. (c) Show that $S^{\perp\perp} = \overline{\langle S \rangle}$, which is the smallest closed subspace of H containing S . (d) Show that for any subspace W of H , $W^{\perp\perp} = \overline{W}$.
55. *Gram–Schmidt Algorithm in a Hilbert Space.* Let $(x_1, x_2, \dots, x_m, \dots)$ be an infinite sequence of linearly independent vectors in H . Give a constructive procedure for computing a sequence $(z_1, z_2, \dots, z_m, \dots)$ of orthonormal vectors in H such that for all $m \geq 1$, (x_1, \dots, x_m) and (z_1, \dots, z_m) span the same subspace. Conclude that every infinite-dimensional Hilbert space has an infinite orthonormal set.
56. (a) Verify that the set L of closed subspaces of H is a complete lattice. (b) Explain why the map $W \mapsto W^\perp$ for $W \in L$ is a lattice anti-isomorphism. (c) Prove or disprove: given $\{W_i : i \in I\} \subseteq L$, the least upper bound of this collection must be the sum of subspaces $\sum_{i \in I} W_i$.
57. (a) Show that if W and Z are closed subspaces of H such that $\langle w, z \rangle = 0$ for all $w \in W$ and $z \in Z$, then $W + Z$ is a closed subspace. (b) Show that every one-dimensional subspace of H is closed. (c) Use (a) and (b) and the existence of orthonormal bases to show that every finite-dimensional subspace of H must be closed.
58. Let X be a set, let $w_k, z_k \in \mathbb{C}$ for each $k \in X$, and let $c \in \mathbb{C}$. (a) Prove $\sum_{k \in X} w_k + \sum_{k \in X} z_k = \sum_{k \in X} (w_k + z_k)$ provided both sums on the left side are finite. (Proceed in three steps, assuming first that all $w_k, z_k \geq 0$, next that all w_k, z_k are real, and then handling the complex case.) (b) Prove

$c \sum_{k \in X} w_k = \sum_{k \in X} (cw_k)$ when the sum on the left is finite. (c) Prove: if $\sum_{k \in X} |w_k| < \infty$, then $\sum_{k \in X} w_k$ is finite and $|\sum_{k \in X} w_k| \leq \sum_{k \in X} |w_k|$. (d) Prove: if $0 \leq w_k \leq z_k$ for each $k \in X$, then $\sum_{k \in X} w_k \leq \sum_{k \in X} z_k$.

59. In §14.10, prove the remaining axioms for the inner product.
60. (a) Prove the polarization identity (14.5). (b) More generally, for any integer $n \geq 3$, prove that for all $x, y \in H$,

$$\sum_{k=0}^{n-1} e^{2\pi i k/n} \|x + e^{2\pi i k/n} y\|^2 = n \langle x, y \rangle.$$

- (c) For $x, y \in H$, evaluate the integral $\int_0^{2\pi} e^{it} \|x + e^{it} y\|^2 dt$.
61. Let V be a complex Banach space in which the norm satisfies the parallelogram law (14.4). Take the polarization identity (14.5) as the definition of $\langle x, y \rangle$ for $x, y \in V$. Show that the axioms for a complex inner product hold, and show that $\langle x, x \rangle = \|x\|^2$, where $\|x\|$ is the given norm on V . (Informally, this says that Hilbert spaces are the Banach spaces where the parallelogram law holds.)
62. Give an example to show that the parallelogram law is not true in \mathbb{C}^n if we use the 1-norm or the sup-norm (see §10.4). Conclude that these normed vector spaces do not have the structure of a Hilbert space.
63. If C and D are closed subsets of H , must $C + D = \{x + y : x \in C, y \in D\}$ be closed in H ?
64. For fixed $z \in H$, define $T_z : H \rightarrow H$ by $T_z(x) = z + x$ for $x \in H$. (a) Prove T_z is an isometry with inverse T_{-z} . (b) Prove T_z maps closed sets to closed sets and convex sets to convex sets.
65. Given a closed subspace W and $x \in H$, we saw that $x = y + z$ for unique $y \in W$ and $z \in W^\perp$. Prove that z is the closest point to x in W^\perp .
66. Given a closed subspace W , define maps $P : H \rightarrow H$ and $Q : H \rightarrow H$ as follows. Given $x \in H$, we know $x = y + z$ for unique $y \in W$ and $z \in W^\perp$. Define $P(x) = y$ and $Q(x) = z$. (a) Prove that P and Q are linear maps such that $P \circ P = P$, $Q \circ Q = Q$, and $P \circ Q = Q \circ P = 0$. (b) Prove that P and Q are continuous maps. (c) Find the image and kernel of P and Q .
67. Let H be the Hilbert space of Lebesgue integrable functions with domain $[0, 2\pi]$. For $n \in \mathbb{Z}$, define $u_n \in H$ by $u_n(t) = e^{int}/\sqrt{2\pi}$ for $t \in [0, 2\pi]$. Prove that $\{u_n : n \in \mathbb{Z}\}$ is an orthonormal set in H . (This orthonormal set, which can be shown to be maximal, plays a central role in the theory of Fourier series.)
68. Let $S = \{u_n : n \in \mathbb{N}^+\}$ be an infinite orthonormal set in H . (a) Compute $\|u_n - u_m\|$ for all $n \neq m$. (b) Prove that a sequence (x_n) with all $x_n \in S$ converges in H iff the sequence is eventually constant, i.e., for some $n_0 \in \mathbb{N}^+$, $x_n = x_{n_0}$ for all $n \geq n_0$. (c) Using (b), show that S is closed and bounded in H , but not compact. (d) Show that any closed set C in H that contains $B(0; \epsilon)$ for some $\epsilon > 0$ is not compact. (Part (d) shows that infinite-dimensional Hilbert spaces are not *locally compact*.)
69. Use Zorn's Lemma (see §16.6) to prove that in any Hilbert space H , there exists a maximal orthonormal set.
70. Let X be an orthonormal subset of H . (a) Prove: if $\|w\|^2 = \sum_{x \in X} |\langle w, x \rangle|^2$ for all $w \in H$, then X is a maximal orthonormal set. (b) Prove: if $\sum_{x \in X} \langle w, x \rangle \overline{\langle z, x \rangle} = \langle w, z \rangle$ for all $w, z \in H$, then X is a maximal orthonormal set.

71. (a) Prove: if $h : X \rightarrow Y$ is a bijection, then there is a Hilbert space isomorphism from $\ell_2(X)$ to $\ell_2(Y)$ sending $f \in \ell_2(X)$ to $f \circ h^{-1} \in \ell_2(Y)$. (b) Prove: if $i : X \rightarrow Y$ is an injection, then i induces a linear isometry mapping $\ell_2(X)$ onto a closed subspace of $\ell_2(Y)$. (c) Prove: for any two Hilbert spaces H_1 and H_2 , H_1 is isometrically isomorphic to a closed subspace of H_2 , or vice versa.
72. Assume V, W, Z are normed vector spaces, $T \in B(V, W)$, and $U \in B(W, Z)$.
- Prove $\|T(y)\| \leq \|T\| \cdot \|y\|$ for all $y \in V$.
 - Prove $\|T\| = \sup\{\|T(x)\|/\|x\| : 0 \neq x \in V\}$.
 - Prove $\|T\| = \inf\{M \in \mathbb{R}^+ : \|T(y)\| \leq M\|y\| \text{ for all } y \in V\}$.
 - Prove $U \circ T \in B(V, Z)$ and $\|U \circ T\| \leq \|U\| \cdot \|T\|$.
73. Given a nonzero $f \in H^*$, let $W = \ker(f)$. Prove W^\perp is one-dimensional.
74. Check that H^* is a Hilbert space with the inner product $\langle R_x, R_y \rangle_{H^*} = \langle y, x \rangle_H$ for $x, y \in H$.
75. Define $E : H \rightarrow H^{**}$ by letting $E(z)$ be the evaluation map $E_z : H^* \rightarrow \mathbb{C}$ given by $E_z(f) = f(z)$ for $f \in H^*$ and $z \in H$. (a) In §14.18, we defined a semi-linear bijective isometry $R_H : H \rightarrow H^*$. Prove that $E = R_{H^*} \circ R_H$. (b) Conclude that E is an isometric isomorphism, so that $H \cong H^{**}$ as Hilbert spaces.
76. Prove: (a) for $S, T \in B(H)$ and $c \in \mathbb{C}$, $(S + T)^* = S^* + T^*$ and $(cS)^* = \bar{c}(S^*)$; (b) for $T \in B(H)$, $\|T^* \circ T\| = \|T\|^2 = \|T \circ T^*\|$ (calculate $\|T(x)\|^2$).
77. (a) Given $T \in B(H)$, show that the *dual map* $T^* : H^* \rightarrow H^*$, given by $T^*(f) = f \circ T$ for $f \in H^*$, is in $B(H^*)$. (b) Prove that dual maps satisfy properties analogous to the properties of adjoint maps listed in §14.19. (c) Let $R : H \rightarrow H^*$ be the semi-isomorphism from §14.18. Prove: for all $T \in B(H)$, $T^* = R^{-1} \circ T^* \circ R$.
78. (a) Prove that the set of self-adjoint operators on H is a closed *real* subspace of $B(H)$. (b) If $S, T \in B(H)$ are self-adjoint, must $S \circ T$ be self-adjoint?
79. (a) Prove: if $T \in B(H)$ satisfies $\langle T(x), x \rangle = 0$ for all $x \in H$, then $T = 0$. (b) Prove: if $S, T \in B(H)$ satisfy $\langle T(x), x \rangle = \langle S(x), x \rangle$ for all $x \in H$, then $S = T$. (c) Prove: $T \in B(H)$ is self-adjoint iff $\langle T(x), x \rangle$ is real for all $x \in H$.
80. Let Z be the set of normal operators in $B(H)$. (a) Is Z closed under addition? scalar multiplication? composition? (b) Is Z a closed set? (c) Prove $T \in Z$ iff $\|T^*(x)\| = \|T(x)\|$ for all $x \in H$. (d) Prove: if $T \in Z$, then $\|T^2\| = \|T\|^2$.
81. Define $T : \ell_2(\mathbb{N}^+) \rightarrow \ell_2(\mathbb{N}^+)$ by setting $T((c_1, c_2, \dots)) = (0, c_1, c_2, \dots)$ for $c_n \in \mathbb{C}$. (a) Prove $T \in B(\ell_2(\mathbb{N}^+))$ and compute $\|T\|$. (b) Find an explicit formula for the map T^* . (c) Show that $T^* \circ T = \text{id}$, but $T \circ T^* \neq \text{id}$. (d) Is T self-adjoint? normal? unitary?
82. Prove that $T \in B(H)$ is unitary iff $T : H \rightarrow H$ is an isometric isomorphism.

Part V

Modules, Independence, and Classification Theorems

This page intentionally left blank

Finitely Generated Commutative Groups

A major goal of abstract algebra is the classification of algebraic structures. An example of such a classification is the theorem of linear algebra stating that every finite-dimensional real vector space V is isomorphic to \mathbb{R}^n for some natural number n . Moreover, the number n is uniquely determined by V and is the dimension of V as a vector space. This classification theorem is useful because, in principle, we can use it to reduce the study of abstract vector spaces such as V to the concrete, familiar vector spaces \mathbb{R}^n .

One could hope for similar classification theorems in the theory of groups. A complete classification of all groups seems far too difficult to ever be achieved. However, much more can be said about special classes of groups, such as simple groups or commutative groups. For example, one can prove that *every finite commutative group is isomorphic to a direct product of cyclic groups, where each cyclic group has size p^e for some prime p and some $e \in \mathbb{N}^+$* . This chapter presents an exposition of a more general classification theorem that describes the structure of all *finitely generated commutative groups*. The techniques used to obtain this structural result for groups can also be applied in a linear-algebraic context to derive canonical forms for matrices and linear transformations. Further abstraction of these arguments eventually leads to a structure theorem classifying finitely generated modules over principal ideal domains, which is one of the keystones of abstract algebra. These topics (modules over PIDs and canonical forms) will be explored in Chapter 18.

To obtain our structure theorem for finitely generated commutative groups, we will first develop the theory of *free commutative groups*. Free commutative groups are the analogues in group theory of the vector spaces that appear so prominently in linear algebra. Thus, our development of free commutative groups will have a distinctively linear-algebraic flavor. In particular, integer-valued matrices will emerge as a key tool for understanding group homomorphisms between two free commutative groups. The reader should constantly bear in mind the analogy between the material presented here and the parallel theory of vector spaces, linear transformations, and matrices.

15.1 Commutative Groups

For convenience, we begin by reviewing some basic facts about commutative groups, which are covered in more detail in Chapter 1. A *commutative group* consists of a set G together with a binary operation on G (denoted by the symbol $+$), subject to the following five conditions. First, for every $g, h \in G$, we must have $g + h \in G$ (closure). Second, for all $g, h, k \in G$, we must have $(g + h) + k = g + (h + k)$ (associativity). Third, for all $g, h \in G$, we must have $g + h = h + g$ (commutativity). Fourth, there must exist an element $0 \in G$ such that $g + 0 = g = 0 + g$ for all $g \in G$ (additive identity). Fifth, for every $g \in G$, there must exist an element $-g \in G$ such that $g + (-g) = 0 = (-g) + g$ (additive inverses). Some familiar examples of commutative groups are the number systems \mathbb{Z} , \mathbb{Q} , \mathbb{R} , and \mathbb{C} , where the group operation is ordinary addition of numbers.

Let H be a subset of a commutative group G . H is called a *subgroup* of G iff $0 \in H$; $h+k \in H$ whenever $h, k \in H$; and $-h \in H$ whenever $h \in H$. Given such a subgroup, we can then form the *quotient group* G/H whose elements are the *cosets* $x+H = \{x+h : h \in H\}$, with x ranging over G . For $x, y \in G$, we have $x+H = y+H$ iff $x-y \in H$. The group operation in G/H is given by

$$(x+H) + (y+H) = (x+y)+H \quad (x, y \in G).$$

See §1.6 for a fuller discussion of quotient groups.

For example, let $G = \mathbb{Z}$, the commutative group of integers under addition. For later work, it will be necessary to know what the subgroups of \mathbb{Z} are. We claim that *every subgroup of \mathbb{Z} has the form $n\mathbb{Z} = \{ni : i \in \mathbb{Z}\}$ for some uniquely determined integer $n \geq 0$* . A routine verification shows that all the sets of the form $n\mathbb{Z}$ are indeed subgroups, no two of which are equal (assuming $n \geq 0$). Now let H be any subgroup of \mathbb{Z} . If H consists of zero alone, then $H = 0\mathbb{Z}$. Otherwise, there must exist a nonzero element $h \in H$. Since $h \in H$ iff $-h \in H$, we see that H contains at least one *positive* integer. Using the least natural number axiom, let n be the smallest positive integer belonging to H . We claim that $n\mathbb{Z} = H$. The inclusion $n\mathbb{Z} \subseteq H$ follows by a quick induction argument, since $n \in H$ and H is closed under addition and inverses. Let us prove the reverse inclusion $H \subseteq n\mathbb{Z}$. Fix $h \in H$. Dividing h by the nonzero integer n , we can write $h = nq+r$ for some integers q, r such that $0 \leq r < n$. Now, since $n\mathbb{Z} \subseteq H$, we have $nq \in H$. It follows that $r = h - (nq) \in H$ as well, since H is closed under addition and inverses. By the minimality of n , this is only possible if $r = 0$. But then $h = nq \in n\mathbb{Z}$. This completes the proof that $H = n\mathbb{Z}$.

For each $n \geq 1$, the quotient group $\mathbb{Z}/n\mathbb{Z}$ gives one way of defining the commutative group of *integers modulo n* . Using the division algorithm in \mathbb{Z} , one can check that $\mathbb{Z}/n\mathbb{Z}$ consists of the n distinct cosets

$$\bar{0} = 0 + n\mathbb{Z}, \bar{1} = 1 + n\mathbb{Z}, \dots, \bar{n-1} = (n-1) + n\mathbb{Z}.$$

Furthermore, this quotient group is isomorphic to the group $\mathbb{Z}_n = \{0, 1, 2, \dots, n-1\}$ with operation \oplus (addition modulo n), as defined in §1.1. Observe that the quotient group $\mathbb{Z}/0\mathbb{Z}$ is isomorphic to \mathbb{Z} itself.

Suppose G_1, \dots, G_k are commutative groups. Consider the product set

$$G_1 \times G_2 \times \dots \times G_k = \{(x_1, x_2, \dots, x_k) : x_i \in G_i \quad (1 \leq i \leq k)\}.$$

This product set becomes a commutative group under the operation

$$(x_1, x_2, \dots, x_k) + (y_1, y_2, \dots, y_k) = (x_1 + y_1, x_2 + y_2, \dots, x_k + y_k) \quad (x_i, y_i \in G_i),$$

as is readily verified. We call this group the *direct product* of the G_i 's. When dealing with commutative groups, the direct product is also written $G_1 \oplus G_2 \oplus \dots \oplus G_k$ and called the (external) *direct sum* of the G_i 's. If every G_i equals the same group G , we may write G^k instead of $G_1 \times \dots \times G_k$. The main goal of this chapter is to prove that every finitely generated commutative group (as defined in the next section) is isomorphic to a product group of the form

$$\mathbb{Z}^b \times \mathbb{Z}_{a_1} \times \mathbb{Z}_{a_2} \times \dots \times \mathbb{Z}_{a_s},$$

where $b \geq 0$, $s \geq 0$, and every a_i is a prime power. Moreover, b , s , and the a_i 's are uniquely determined by G .

15.2 Generating Sets

Suppose G is a commutative group, $x \in G$, and n is an integer. If $n > 0$, define $nx = x + x + \cdots + x$ (n summands). If $n = 0$, define $nx = 0$. If $n < 0$, define $nx = -x + -x + \cdots + -x$ ($|n|$ summands). Group elements of the form nx are called *multiples* of x . For all $x, y \in G$ and all $m, n \in \mathbb{Z}$, the following rules hold:

$$(m+n)x = (mx) + (nx); \quad m(nx) = (mn)x; \quad 1x = x; \quad -(mx) = (-m)x; \quad (15.1)$$

$$m(x+y) = (mx) + (my).$$

The first four rules are none other than the “laws of exponents” (which are valid for arbitrary groups), translated from multiplicative to additive notation. The final rule is only valid because the group G is commutative. Informally, this rule holds when $m > 0$ since

$$\begin{aligned} m(x+y) &= (x+y) + (x+y) + \cdots + (x+y) \quad (m \text{ summands}) \\ &= (\underbrace{x+x+\cdots+x}_m) + (\underbrace{y+y+\cdots+y}_m) \quad (\text{since } G \text{ is commutative}) \\ &= (mx) + (my). \end{aligned}$$

A more formal verification of this rule (see Exercise 6) uses induction to establish its validity for $m \geq 0$, followed by a separate argument for negative m .

Next, suppose v_1, \dots, v_k are elements of a commutative group G . A \mathbb{Z} -linear combination of these elements is an element of the form

$$c_1v_1 + c_2v_2 + \cdots + c_kv_k,$$

where each c_i is an integer. The set H of all \mathbb{Z} -linear combinations of v_1, \dots, v_k is a subgroup of G , denoted by $\langle v_1, v_2, \dots, v_k \rangle$. This follows from the rules given above and the fact that addition is commutative in G . For instance, H is closed under addition since (for $c_i, d_i \in \mathbb{Z}$)

$$\sum_{i=1}^k c_i v_i + \sum_{i=1}^k d_i v_i = \sum_{i=1}^k (c_i v_i + d_i v_i) = \sum_{i=1}^k (c_i + d_i) v_i.$$

We call $\langle v_1, v_2, \dots, v_k \rangle$ the *subgroup of G generated by v_1, \dots, v_k* . If there exist finitely many elements v_1, \dots, v_k that generate G , then G is called a *finitely generated* commutative group. If G can be generated by a single element v_1 , then G is called a *cyclic* group. Given any commutative group G and any element $x \in G$, the set of multiples $\langle x \rangle = \{nx : n \in \mathbb{Z}\}$ is a cyclic subgroup of G . We sometimes denote this subgroup by the symbol $\mathbb{Z}x$.

Every finite commutative group G is finitely generated, since we can take all the elements of G as a generating set. Of course, there are likely to be other generating sets that are much smaller. The groups \mathbb{Z} and $\mathbb{Z}/n\mathbb{Z}$ (for $n \geq 1$) are finitely generated groups. Indeed, they are cyclic groups since $\mathbb{Z} = \langle 1 \rangle$ and $\mathbb{Z}/n\mathbb{Z} = \langle 1 + n\mathbb{Z} \rangle$. Generators of cyclic groups are not unique; for instance, $\mathbb{Z} = \langle 1 \rangle = \langle -1 \rangle$. Also, one can show that for prime p , \mathbb{Z}_p is generated by any of its elements other than 0 (Exercise 10).

Let G_1, \dots, G_k be finitely generated commutative groups, say $G_i = \langle v_{i,1}, \dots, v_{i,n_i} \rangle$ for some $n_i \in \mathbb{N}$ and $v_{i,j} \in G_i$. Given the product group $G = G_1 \times \cdots \times G_k$, we can associate to each $v_{i,j}$ the k -tuple $(0, 0, \dots, v_{i,j}, \dots, 0)$, where $v_{i,j}$ appears in position i . One can check that G is generated by these k -tuples, hence is finitely generated. For example, since $\mathbb{Z} = \langle 1 \rangle$, \mathbb{Z}^k is generated by the k elements $e_i = (0, 0, \dots, 1, \dots, 0)$, where the 1 occurs in position i . Note that e_i is the image in \mathbb{Z}^k of the generator 1 of the i 'th copy of \mathbb{Z} .

Not every commutative group is finitely generated. For example, consider the group \mathbb{Q} of rational numbers. To see that \mathbb{Q} (under addition) cannot be finitely generated, argue by contradiction. Let $v_1, \dots, v_k \in \mathbb{Q}$ be a finite generating set. Write $v_i = a_i/b_i$ for some integers a_i, b_i with $b_i > 0$. By finding a common denominator, we can change notation so that $v_i = c_i/d$ for some integers c_i and d with $d > 0$. (For instance, let $d = b_1 b_2 \cdots b_k$ and $c_i = (da_i/b_i) \in \mathbb{Z}$.) Now consider an arbitrary \mathbb{Z} -linear combination of the v_i 's:

$$n_1 v_1 + \cdots + n_k v_k = \frac{n_1 c_1 + \cdots + n_k c_k}{d} \quad (n_i \in \mathbb{Z}).$$

Any such number is either 0 or has absolute value at least $1/d$. But then the rational number $1/(2d)$ cannot be expressed in this form, because its absolute value is too small.

Although this chapter is concerned mainly with finitely generated commutative groups, it is possible to define the notion of an infinite generating set. Let S be an arbitrary subset (possibly infinite) of a commutative group G . By definition, a \mathbb{Z} -linear combination of elements of S is any \mathbb{Z} -linear combination of any finite subset S' of S . We often write $\sum_{x \in S} n_x x$ to denote such a linear combination, understanding that all but a finite number of the coefficients $n_x \in \mathbb{Z}$ must be zero. One can check, as before, that the set of \mathbb{Z} -linear combinations of elements of S is a subgroup of G , denoted $\langle S \rangle$. We say S generates G iff $G = \langle S \rangle$. For example, using the fundamental theorem of arithmetic, one can verify that the set of rational numbers of the form $1/p^e$ (for p prime and $e \geq 1$) generates the commutative group \mathbb{Q} .

15.3 \mathbb{Z} -Independence and \mathbb{Z} -Bases

A generating set for a commutative group resembles a spanning set for a vector space. The only difference is that the “scalars” are integers rather than field elements. The analogy to linear algebra suggests the following definitions. Let (v_1, \dots, v_k) be a finite list of elements in a commutative group G . We call this list \mathbb{Z} -linearly independent (or \mathbb{Z} -independent) iff for all integers c_1, \dots, c_k ,

$$\text{if } c_1 v_1 + c_2 v_2 + \cdots + c_k v_k = 0_G, \text{ then } c_1 = c_2 = \cdots = c_k = 0.$$

This means that no linear combination of the v_i 's can produce 0_G except the one where all coefficients are zero. If S is a subset of G (possibly infinite), we say that S is \mathbb{Z} -linearly independent iff every finite nonempty list of distinct elements of S is \mathbb{Z} -independent.

An ordered \mathbb{Z} -basis of G is a list $B = (v_1, \dots, v_k)$ of elements $v_i \in G$ such that B is \mathbb{Z} -independent and $G = \langle v_1, \dots, v_k \rangle$. A subset S of G (finite or not) is called a \mathbb{Z} -basis of G iff $G = \langle S \rangle$ and S is \mathbb{Z} -linearly independent. If G has a basis, then G is called a free commutative group. If G has a k -element basis for some $k \in \mathbb{N}$, G is a finitely generated free commutative group with dimension (or rank) k . Later, we will prove that the dimension of such a group is uniquely determined.

Not every commutative group is free. For instance, suppose G is a finite group with more than one element. Consider any list (v_1, \dots, v_k) of elements of G . From group theory, there exists an integer $n > 0$ with $nv_1 = 0$; for instance, $n = |G|$ has this property (see Exercise 8 in Chapter 1 for a proof). Then the relation $nv_1 + 0v_2 + \cdots + 0v_k = 0$ shows that the given list must be \mathbb{Z} -linearly dependent. On the other hand, if $G = \{0\}$ is the one-element group, then the empty set is a \mathbb{Z} -basis of G . This follows from the convention $\langle \emptyset \rangle = \{0\}$ and the fact that the empty set is \mathbb{Z} -independent, which is a logical consequence of the definition.

For more interesting examples of bases, consider the group \mathbb{Z}^k of all k -tuples of integers, with group operation

$$(a_1, a_2, \dots, a_k) + (b_1, b_2, \dots, b_k) = (a_1 + b_1, a_2 + b_2, \dots, a_k + b_k) \quad (a_i, b_i \in \mathbb{Z}).$$

For $1 \leq i \leq k$, let $e_i = (0, \dots, 0, 1, 0, \dots, 0)$, where the 1 occurs in position i . Let us verify explicitly that the list $B = (e_1, e_2, \dots, e_k)$ is an ordered basis of \mathbb{Z}^k , called the *standard ordered basis of \mathbb{Z}^k* . First, B spans \mathbb{Z}^k , because for any $(a_1, \dots, a_k) \in \mathbb{Z}^k$, we have

$$(a_1, \dots, a_k) = \sum_{i=1}^k a_i e_i.$$

Second, B is \mathbb{Z} -independent, because the relation $\sum_{i=1}^k b_i e_i = 0$ (with $b_i \in \mathbb{Z}$) means that $(b_1, \dots, b_k) = (0, \dots, 0)$, which implies that all b_i 's are zero by equating corresponding components. Thus, we have proved that \mathbb{Z}^k is a k -dimensional free commutative group.

15.4 Elementary Operations on \mathbb{Z} -Bases

As in the case of vector spaces, most free commutative groups have many different bases. There are three *elementary operations* that allow us to produce new ordered bases from old ones. We now define how each of these operations affects a given ordered \mathbb{Z} -basis $B = (v_1, \dots, v_k)$ of a finitely generated free commutative group G .

(B1) For any $i \neq j$, we can interchange the position of v_i and v_j in the ordered list B . Using commutativity, we see (Exercise 20) that the new ordered list still generates G and is still \mathbb{Z} -independent, so it is a \mathbb{Z} -basis of G .

(B2) For any i , we can replace v_i in B by $-v_i$. Since $c_i v_i = (-c_i)(-v_i)$, one readily checks (Exercise 21) that the new ordered list is still a \mathbb{Z} -basis of G .

(B3) For any $i \neq j$ and any integer c , we can replace v_i in B by $w_i = v_i + cv_j$. Let us check that this gives another \mathbb{Z} -basis of G . In light of (B1), it suffices to consider the case $i = 1$ and $j = 2$. Write $B = (v_1, v_2, \dots, v_k)$ and $B' = (v_1 + cv_2, v_2, \dots, v_k)$. First, does B' generate G ? Given any $g \in G$, we can write

$$g = n_1 v_1 + n_2 v_2 + \cdots + n_k v_k$$

for some $n_i \in \mathbb{Z}$, because B is known to generate G . Manipulating this expression gives

$$g = n_1(v_1 + cv_2) + (n_2 - cn_1)v_2 + n_3 v_3 + \cdots + n_k v_k$$

where all coefficients are integers. Thus, g is a linear combination of the elements in the list B' . Second, is B' \mathbb{Z} -independent? Assume that

$$d_1(v_1 + cv_2) + d_2 v_2 + \cdots + d_k v_k = 0$$

for fixed $d_i \in \mathbb{Z}$; we must prove every $d_i = 0$. Rewriting the assumption gives

$$d_1 v_1 + (cd_1 + d_2) v_2 + d_3 v_3 + \cdots + d_k v_k = 0.$$

By the known \mathbb{Z} -independence of B , we conclude that $d_1 = cd_1 + d_2 = d_3 = \cdots = d_k = 0$. Then $d_2 = (cd_1 + d_2) - cd_1 = 0 - c0 = 0$. So all d_i 's are zero.

For example, starting with the standard ordered basis (e_1, e_2, e_3) of \mathbb{Z}^3 , we can apply a sequence of elementary operations to produce new ordered bases of this group. As a specific illustration, one can check that the list $(e_2 - 5e_3, -e_3, e_1 + 2e_2 + 3e_3) = ((0, 1, -5), (0, 0, -1), (1, 2, 3))$ can be obtained from (e_1, e_2, e_3) by appropriate elementary operations. Hence, this list is an ordered \mathbb{Z} -basis of \mathbb{Z}^3 .

We remark that similar elementary operations can be applied to ordered bases of vector spaces over a field F . In operation (B3), we replace the integer c by a scalar $c \in F$. In operation (B2), we may now select any nonzero scalar $c \in F$ and replace v_i by cv_i . One can check that the corresponding operation for commutative groups only produces a \mathbb{Z} -basis when $c = \pm 1$. Essentially, the reason is that $+1$ and -1 are the only integers whose *multiplicative* inverses are again integers.

15.5 Coordinates and \mathbb{Z} -Linear Maps

The next few sections derive some fundamental facts about bases and free commutative groups. We focus attention on the case of finitely generated groups, but all our proofs extend without difficulty to the case of groups with infinite bases.

Let G be a free commutative group with ordered basis $B = (v_1, \dots, v_k)$. We first prove: *for every $g \in G$, there exist unique integers n_1, \dots, n_k such that $g = n_1v_1 + \dots + n_kv_k$.* We sometimes call (n_1, \dots, n_k) the *coordinates of g relative to B* . The *existence* of the integers n_i follows immediately from the fact that the v_i 's generate G . As for *uniqueness*, suppose $g \in G$ and

$$n_1v_1 + \dots + n_kv_k = g = m_1v_1 + \dots + m_kv_k$$

for some $n_i, m_i \in \mathbb{Z}$. Rearranging and using the rules in (15.1), this equation becomes

$$(n_1 - m_1)v_1 + \dots + (n_k - m_k)v_k = 0.$$

By \mathbb{Z} -independence of B , it follows that $n_i - m_i = 0$ for all i . Thus, $n_i = m_i$ for all i , proving that the coordinates n_i are uniquely determined by g and B .

To continue the development of free commutative groups, we need the idea of a \mathbb{Z} -linear map. By definition, a *\mathbb{Z} -linear map* is a group homomorphism $T : G \rightarrow H$ between two commutative groups G and H . To say that T is a homomorphism means that $T(x+y) = T(x) + T(y)$ for all $x, y \in G$. It follows by a routine induction argument that $T(nx) = nT(x)$ for all $n \in \mathbb{Z}$ and all $x \in G$, which explains the terminology “ \mathbb{Z} -linear.” It also follows by induction that a \mathbb{Z} -linear map T preserves \mathbb{Z} -linear combinations, i.e.,

$$T\left(\sum_{i=1}^k n_i v_i\right) = \sum_{i=1}^k n_i T(v_i) \quad (k \in \mathbb{N}^+, n_i \in \mathbb{Z}, v_i \in G).$$

Note the close analogy to linear transformations between vector spaces, which satisfy similar identities. As in the case of linear transformations, the *kernel* of a \mathbb{Z} -linear map $T : G \rightarrow H$ is $\ker(T) = \{x \in G : T(x) = 0_H\}$, and the *image* of T is $\text{img}(T) = \{T(x) : x \in G\}$. The \mathbb{Z} -linear map T is injective iff $\ker(T) = \{0_G\}$; T is surjective iff $\text{img}(T) = H$. The *fundamental homomorphism theorem* states that any \mathbb{Z} -linear map $T : G \rightarrow H$ with kernel K and image I induces a \mathbb{Z} -linear isomorphism $T' : G/K \rightarrow I$ given by $T'(x+K) = T(x)$ for all $x \in G$. See §1.7 for a proof of this theorem.

We pause to establish two basic facts about \mathbb{Z} -linear maps. Let $T : G \rightarrow H$ be a \mathbb{Z} -linear map between two commutative groups. Assume H is generated by elements w_1, \dots, w_m . The

first fact says that *the image of T is all of H iff every generator w_j lies in the image of T .* The forward implication is immediate; conversely, suppose each $w_j = T(x_j)$ for some $x_j \in G$. Given $h \in H$, h can be written (not necessarily uniquely) in the form $h = n_1w_1 + \cdots + n_mw_m$ where $n_i \in \mathbb{Z}$. Choosing $x = n_1x_1 + \cdots + n_mx_m \in G$, \mathbb{Z} -linearity implies that $T(x) = h$. To state the second fact, let $S : G \rightarrow H$ be another \mathbb{Z} -linear map, and assume G is generated by elements v_1, \dots, v_k . Then $T = S$ iff $T(v_i) = S(v_i)$ for all i with $1 \leq i \leq k$. In other words, *two \mathbb{Z} -linear maps are equal iff they agree on a generating set for the domain.* The forward implication is immediate. To prove the converse, suppose $T(v_i) = S(v_i)$ for all i , and let $g \in G$. We can write $g = p_1v_1 + \cdots + p_kv_k$ for some integers p_i . Using \mathbb{Z} -linearity and the hypothesis on T and S , we see that

$$T(g) = \sum_{i=1}^k p_i T(v_i) = \sum_{i=1}^k p_i S(v_i) = S(g).$$

15.6 UMP for Free Commutative Groups

Now we are ready to prove a theorem of a fundamental nature regarding free commutative groups. This theorem is called the *universal mapping property (UMP)* for free commutative groups. Suppose $X = \{v_1, \dots, v_k\}$ is a basis of a free commutative group G . For every commutative group H and every function $f : X \rightarrow H$, there exists a unique \mathbb{Z} -linear map $T_f : G \rightarrow H$ such that $T_f(v_i) = f(v_i)$ for all $v_i \in X$. We call T_f the \mathbb{Z} -linear extension of f from X to G .

Let us prove the existence of T_f first. To define the value of T_f at some given $g \in G$, first write g in the form

$$g = n_1v_1 + \cdots + n_kv_k \quad (n_j \in \mathbb{Z}).$$

We have already proved that the n_j 's appearing in this expression are uniquely determined by g . Therefore, the formula

$$T_f(g) = n_1f(v_1) + \cdots + n_kf(v_k)$$

gives a well-defined element of H for each $g \in G$. We need only check that the function T_f is a group homomorphism extending f . First, is $T_f(g_1 + g_2) = T_f(g_1) + T_f(g_2)$ for all $g_1, g_2 \in G$? To answer this, first write $g_1 = \sum_i n_i v_i$ and $g_2 = \sum_i m_i v_i$ where $n_i, m_i \in \mathbb{Z}$. By the definition of T_f , we then have

$$T_f(g_1) = n_1f(v_1) + \cdots + n_kf(v_k);$$

$$T_f(g_2) = m_1f(v_1) + \cdots + m_kf(v_k).$$

On the other hand, note that $g_1 + g_2 = \sum_i (n_i + m_i)v_i$ is the unique expression for $g_1 + g_2$ as a linear combination of the v_i 's. So, applying the definition of T_f gives

$$T_f(g_1 + g_2) = (n_1 + m_1)f(v_1) + \cdots + (n_k + m_k)f(v_k).$$

Comparing to the previous formulas, it is now evident that $T_f(g_1 + g_2) = T_f(g_1) + T_f(g_2)$. Next, does T_f extend f ? To answer this, let j be a fixed index. Observe that $v_j = 0v_1 + \cdots + 1v_j + \cdots + 0v_k$ is the unique expansion of v_j in terms of the v_i 's. Therefore, by definition,

$$T_f(v_j) = 0f(v_1) + \cdots + 1f(v_j) + \cdots + 0f(v_k) = f(v_j).$$

This completes the existence proof.

Uniqueness of T_f is established as follows. Suppose $S : G \rightarrow H$ is any \mathbb{Z} -linear map extending f . Then $T_f(v_j) = f(v_j) = S(v_j)$ for all j , so T_f and S agree on a generating set for G . So $T_f = S$. This completes the proof of the UMP. One can verify that the UMP holds for all free commutative groups (even those that are not finitely generated).

We can use the UMP to derive some further facts about free commutative groups. For instance, *every k -dimensional free commutative group G is isomorphic to \mathbb{Z}^k* . More precisely, let $B = (v_1, \dots, v_k)$ be an ordered basis of G . Define a function $f : \{v_1, \dots, v_k\} \rightarrow \mathbb{Z}^k$ by setting $f(v_i) = e_i$, where e_i is the k -tuple with a 1 in position i and zeroes elsewhere. Using the UMP, we obtain a unique \mathbb{Z} -linear map $T : G \rightarrow \mathbb{Z}^k$ extending f . Now, T is surjective because the image of T contains the generating set $\{e_1, \dots, e_k\}$ of \mathbb{Z}^k . To see that T is injective, recall the formula defining T :

$$T(n_1 v_1 + \cdots + n_k v_k) = n_1 f(v_1) + \cdots + n_k f(v_k) = n_1 e_1 + \cdots + n_k e_k \quad (n_i \in \mathbb{Z}).$$

If T sends $x = n_1 v_1 + \cdots + n_k v_k \in G$ to zero, then $n_1 = \cdots = n_k = 0$ by the \mathbb{Z} -independence of the e_i 's. Hence, $x = 0$, proving that T has kernel $\{0\}$. It follows that T is a bijective \mathbb{Z} -linear map, which defines an isomorphism from G to \mathbb{Z}^k .

15.7 Quotient Groups of Free Commutative Groups

Our goal in this chapter is to classify all finitely generated commutative groups. More precisely, part of our goal is to show that every finitely generated commutative group is isomorphic to a direct sum of cyclic groups (namely \mathbb{Z} or $\mathbb{Z}/n\mathbb{Z} \cong \mathbb{Z}_n$). We achieved part of this goal at the end of the last section, by showing that every finitely generated *free* commutative group is isomorphic to one of the groups \mathbb{Z}^k . This is analogous to the linear algebra theorem stating that every finite-dimensional vector space over a field F is isomorphic to F^k for some $k \geq 0$. However, in the case of commutative groups, there is more work to do since not all commutative groups are free. (We must also eventually address the question of the *uniqueness* of k .)

Next, we continue to progress towards our goal by proving that *every finitely generated commutative group is isomorphic to a quotient group F/P , for some finitely generated free commutative group F and some subgroup P* . (This statement remains true, with the same proof, if the two occurrences of “finitely generated” are deleted.) Let H be a commutative group generated by w_1, \dots, w_k . Let \mathbb{Z}^k be the free commutative group with its standard ordered basis (e_1, \dots, e_k) . Define a map $f : \{e_1, \dots, e_k\} \rightarrow H$ by setting $f(e_i) = w_i$ for $1 \leq i \leq k$. Next, use the UMP to obtain a \mathbb{Z} -linear extension $T : \mathbb{Z}^k \rightarrow H$. The image of T is all of H , since all the generators w_j of H lie in the image of T . Let P be the kernel of T . Applying the fundamental homomorphism theorem to T , we see that T induces a group isomorphism $T' : \mathbb{Z}^k / P \rightarrow H$. Thus H is isomorphic to a quotient group of a free commutative group whose dimension is the same as the size of the given generating set for H .

To see why this result will help us in the classification of finitely generated commutative groups, consider the special case where the subgroup P of \mathbb{Z}^k has the particular form

$$P = n_1 \mathbb{Z} \times n_2 \mathbb{Z} \times \cdots \times n_k \mathbb{Z}, \tag{15.2}$$

for certain $n_1, \dots, n_k \in \mathbb{N}$. Define a map $S : \mathbb{Z}^k \rightarrow (\mathbb{Z}/n_1\mathbb{Z}) \times \cdots \times (\mathbb{Z}/n_k\mathbb{Z})$ by setting $S((a_1, \dots, a_k)) = (a_1 + n_1 \mathbb{Z}, \dots, a_k + n_k \mathbb{Z})$ for $a_i \in \mathbb{Z}$. One readily checks that S is a

surjective group homomorphism with kernel P . Hence, by the fundamental homomorphism theorem, we obtain an isomorphism

$$H \cong \mathbb{Z}^k / P \cong \mathbb{Z}/n_1\mathbb{Z} \times \mathbb{Z}/n_2\mathbb{Z} \times \cdots \times \mathbb{Z}/n_k\mathbb{Z} \cong \mathbb{Z}_{n_1} \times \mathbb{Z}_{n_2} \times \cdots \times \mathbb{Z}_{n_k}.$$

Thus, provided P is a subgroup of the form (15.2), we have succeeded in writing H as a direct sum of cyclic groups.

Unfortunately, when $k > 1$, not every subgroup of \mathbb{Z}^k has the form given in (15.2). For example, when $k = 2$, consider the subgroup $P = \{(t, 2t) : t \in \mathbb{Z}\} \subseteq \mathbb{Z} \times \mathbb{Z}$. One sees that the set P is not of the form $A \times B$ for any choice of $A, B \subseteq \mathbb{Z}$. In general, the subgroup structure of the free commutative groups \mathbb{Z}^k is rich and subtle. However, the time invested in studying these subgroups will eventually lead us to our goal of classifying finitely generated commutative groups.

Before beginning a closer study of the subgroups of \mathbb{Z}^k , let us revisit the analogy to vector spaces to gain some intuition. The vector-space analogue of the statement “not every P has the form (15.2)” is the statement “not every subspace of \mathbb{R}^k is spanned by a subset of the standard basis vectors e_i .” The vector-space analogue of the subgroup P considered in the last paragraph is the subspace $\{(t, 2t) : t \in \mathbb{R}\}$ of \mathbb{R}^2 . This one-dimensional subspace of \mathbb{R}^2 is not spanned by either of the vectors $(1, 0)$ or $(0, 1)$. However, it is spanned by the vector $(1, 2)$. We know from linear algebra that the one-element set $\{(1, 2)\}$ can be extended to a basis of \mathbb{R}^2 . This suggests the possibility of *changing the basis of \mathbb{Z}^k* to force the subgroup P to assume the nice form given in (15.2). We will pursue this idea in the next few sections.

15.8 Subgroups of Free Commutative Groups

To start implementing the agenda outlined at the end of the last section, we need to prove a key technical point about subgroups of free commutative groups. We will show that *if G is a k -dimensional free commutative group and H is a subgroup of G , then H is also a free commutative group with dimension at most k .* Since we know G is isomorphic to \mathbb{Z}^k , we will prove this result for $G = \mathbb{Z}^k$ without loss of generality. The proof uses induction on k .

The case $k = 0$ is immediate. Suppose $k = 1$. In §15.1, we saw that every subgroup of \mathbb{Z}^1 has the form $n\mathbb{Z}$ for some integer $n \geq 0$. If $n = 0$, then $0\mathbb{Z}$ is a zero-dimensional free commutative group with an empty basis. If $n > 0$, then one can check that $n\mathbb{Z} \cong \mathbb{Z}$ is a one-dimensional free commutative group with a one-element basis $\{n\}$. (We observe in passing that, for $n > 1$, this basis of the subgroup $n\mathbb{Z}$ cannot be extended to a basis of \mathbb{Z} . This reveals one notable difference between free commutative groups and vector spaces.)

For the induction step, suppose $k > 1$ and the theorem is already known for all free commutative groups of dimension less than k . Let H be a fixed subgroup of \mathbb{Z}^k . To apply our induction hypothesis, we need a subgroup of \mathbb{Z}^{k-1} . To obtain such a subgroup, let $H' = H \cap (\mathbb{Z}^{k-1} \times \{0\})$ be the set of all elements of H with last coordinate zero. Note that H' is a subgroup of $\mathbb{Z}^{k-1} \times \{0\}$, which is a $(k-1)$ -dimensional free commutative group isomorphic to \mathbb{Z}^{k-1} . Applying the induction hypothesis, we conclude that H' is free and has some ordered basis (v_1, \dots, v_{m-1}) where $m-1 \leq k-1$. We must somehow pass from this basis to a basis of the full subgroup H .

Toward this end, consider the projection map $P : \mathbb{Z}^k \rightarrow \mathbb{Z}$ given by $P((a_1, \dots, a_k)) = a_k$ for $a_i \in \mathbb{Z}$. It is routine to check that P is a group homomorphism. Therefore, $P[H] = \{P(h) : h \in H\}$ is a subgroup of \mathbb{Z} . By our earlier classification of the subgroups of

\mathbb{Z} , we know there is some integer $q \geq 0$ such that $P[H] = q\mathbb{Z}$. Let v_m be any fixed element of H such that $P(v_m) = q$. So $v_m = (a_1, \dots, a_{k-1}, q) \in H$ for some integers a_i .

Now consider two cases. First, suppose $q = 0$. Then $P[H] = \{0\}$, which means that every element of H has last coordinate zero. But then $H = H'$, and we already know that H' is a free commutative group of dimension $m - 1 \leq k - 1 < k$.

The second case is that $q > 0$. We claim that $X = (v_1, \dots, v_m)$ is an ordered basis of H , so that H is free with dimension $m \leq k$. To prove the claim, we first check the \mathbb{Z} -independence of X . Suppose $c_1v_1 + \dots + c_mv_m = 0$ for some integers c_i . Apply the \mathbb{Z} -linear projection map P to this relation to obtain $c_1P(v_1) + \dots + c_{m-1}P(v_{m-1}) + c_mP(v_m) = P(0) = 0$. For $i < m$, $P(v_i) = 0$ since $v_i \in H'$. On the other hand, $P(v_m) = q$ by choice of v_m . So the relation reduces to $c_mq = 0$, which implies $c_m = 0$ because $q > 0$ and \mathbb{Z} has no zero divisors. However, once we know that $c_m = 0$, the original relation becomes $c_1v_1 + \dots + c_{m-1}v_{m-1} = 0$. We can now conclude that $c_1 = \dots = c_{m-1} = 0$ because v_1, \dots, v_{m-1} are already known to be \mathbb{Z} -linearly independent. To finish the proof, we must show that X generates H . Fix $h = (b_1, \dots, b_k) \in H$, where $b_i \in \mathbb{Z}$. We have $P(h) = b_k \in P[H] = q\mathbb{Z}$, so $b_k = tq$ for some integer t . Note that

$$h - tv_m = (b_1, \dots, b_k) - t(a_1, \dots, a_{k-1}, q) = (b_1 - ta_1, \dots, b_{k-1} - ta_{k-1}, 0),$$

so $h - tv_m \in H'$. Since (v_1, \dots, v_{m-1}) is known to generate H' , we have $h - tv_m = d_1v_1 + \dots + d_{m-1}v_{m-1}$ for certain integers d_i . Then

$$h = d_1v_1 + \dots + d_{m-1}v_{m-1} + tv_m,$$

so that h is a \mathbb{Z} -linear combination of v_1, \dots, v_m .

We should mention that every subgroup of an infinite-dimensional free commutative group is also free, but we will not prove this fact here.

15.9 \mathbb{Z} -Linear Maps and Integer Matrices

We have not yet finished our analysis of the subgroups of \mathbb{Z}^k . However, to complete our work in this area, we must first develop some machinery for understanding general \mathbb{Z} -linear maps. Eventually, we will apply this material to the \mathbb{Z} -linear inclusion map of a subgroup into \mathbb{Z}^k .

Let us consider the following setup. Suppose $T : G \rightarrow H$ is a \mathbb{Z} -linear map between two finitely generated free commutative groups. Let $X = (v_1, \dots, v_n)$ be an ordered basis of G , and let $Y = (w_1, \dots, w_m)$ be an ordered basis of H . We first show how to use the ordered bases X and Y to represent T by an $m \times n$ matrix with integer entries. The following construction is exactly analogous to the procedure used in linear algebra to associate a matrix of scalars to a given linear transformation between two vector spaces (see Chapter 6).

We first remark that the \mathbb{Z} -linear map T is completely determined by its effect on the generators v_1, \dots, v_n of G . Thus, to specify T , we need only record the n elements $T(v_1), \dots, T(v_n)$. Next, for each j , we know that $T(v_j)$ can be expressed *uniquely* as a \mathbb{Z} -linear combination of w_1, \dots, w_m . In other words, for all j with $1 \leq j \leq n$, we can write

$$T(v_j) = \sum_{i=1}^m a_{ij}w_i \tag{15.3}$$

for certain uniquely determined integers a_{ij} . The $m \times n$ matrix $A = [a_{ij}]$ is called the *matrix*

of T relative to the bases X and Y . As long as X and Y are fixed and known, the passage from T to A is completely reversible. In other words, we could start with the matrix A and use (15.3) as the definition of T on the generators v_j . By the UMP for G , this definition extends uniquely to a \mathbb{Z} -linear map from G to H .

For example, consider the matrix

$$A = \begin{bmatrix} 7 & 2 & -1 & 5 \\ 0 & 4 & 0 & 2 \\ 3 & 3 & 1 & 0 \end{bmatrix}.$$

Take $G = \mathbb{Z}^4$, $H = \mathbb{Z}^3$, $v_1 = (1, 0, 0, 0)$, $v_2 = (0, 1, 0, 0)$, $v_3 = (0, 0, 1, 0)$, $v_4 = (0, 0, 0, 1)$, $w_1 = (1, 0, 0)$, $w_2 = (0, 1, 0)$, and $w_3 = (0, 0, 1)$, so that X and Y are the standard ordered bases of \mathbb{Z}^4 and \mathbb{Z}^3 , respectively. Given this data, we obtain a \mathbb{Z} -linear map $T : \mathbb{Z}^4 \rightarrow \mathbb{Z}^3$ defined on basis elements by

$$\begin{aligned} T(v_1) &= 7w_1 + 0w_2 + 3w_3 = (7, 0, 3); \\ T(v_2) &= 2w_1 + 4w_2 + 3w_3 = (2, 4, 3); \\ T(v_3) &= -1w_1 + 0w_2 + 1w_3 = (-1, 0, 1); \\ T(v_4) &= 5w_1 + 2w_2 + 0w_3 = (5, 2, 0). \end{aligned}$$

Note that the j 'th column of A provides the coordinates of the image of the j 'th input basis element relative to the given output basis. We can use \mathbb{Z} -linearity to compute explicitly the image of an arbitrary element $(a, b, c, d) \in \mathbb{Z}^4$. In detail, for any $a, b, c, d \in \mathbb{Z}$,

$$\begin{aligned} T((a, b, c, d)) &= T(av_1 + bv_2 + cv_3 + dv_4) \\ &= aT(v_1) + bT(v_2) + cT(v_3) + dT(v_4) \\ &= a(7, 0, 3) + b(2, 4, 3) + c(-1, 0, 1) + d(5, 2, 0) \\ &= (7a + 2b - c + 5d, 4b + 2d, 3a + 3b + c). \end{aligned}$$

The same answer can be found by computing the following matrix-vector product:

$$\begin{bmatrix} 7 & 2 & -1 & 5 \\ 0 & 4 & 0 & 2 \\ 3 & 3 & 1 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} 7a + 2b - c + 5d \\ 4b + 2d \\ 3a + 3b + c \end{bmatrix}.$$

It is now possible to interpret matrix addition and matrix multiplication in terms of \mathbb{Z} -linear maps. Returning to the general setup, suppose $T, U : G \rightarrow H$ are two \mathbb{Z} -linear maps. Let $A = [a_{ij}]$ and $B = [b_{ij}]$ be the integer matrices representing T and U relative to the ordered bases X and Y . Then, by definition, we must have

$$T(v_j) = \sum_{i=1}^m a_{ij} w_i \quad (1 \leq j \leq n);$$

$$U(v_j) = \sum_{i=1}^m b_{ij} w_i \quad (1 \leq j \leq n).$$

The *sum* of the maps T and U is the map $T+U : G \rightarrow H$ defined by $(T+U)(x) = T(x)+U(x)$ for all $x \in G$; we see at once that $T+U$ is \mathbb{Z} -linear. What is the matrix of $T+U$ relative to X and Y ? To find the j 'th column of this matrix, we must find the coordinates of $(T+U)(v_j)$ relative to Y . We find that

$$(T+U)(v_j) = T(v_j) + U(v_j) = \sum_{i=1}^m a_{ij} w_i + \sum_{i=1}^m b_{ij} w_i = \sum_{i=1}^m (a_{ij} + b_{ij}) w_i.$$

Thus, $T+U$ is represented by the $m \times n$ matrix $[a_{ij} + b_{ij}] = A+B$. This shows that *addition of matrices corresponds to addition of \mathbb{Z} -linear maps*.

Products are a bit more subtle. Take T and A as above, and suppose $V : H \rightarrow K$ is a \mathbb{Z} -linear map from H into a third free commutative group K with ordered basis $Z = (z_1, \dots, z_p)$. There are unique integers c_{ik} such that

$$V(w_k) = \sum_{i=1}^p c_{ik} z_i \quad (1 \leq k \leq m);$$

then $C = [c_{ik}]$ is the $p \times m$ matrix of V relative to the ordered bases Y and Z . Now, we know that the composite map $V \circ T : G \rightarrow K$ is \mathbb{Z} -linear; what is the matrix of this map relative to the ordered bases X and Z ? As before, we discover the answer by applying this map to a typical basis vector $v_j \in X$. Using \mathbb{Z} -linearity of the maps, commutativity of addition, and the distributive law, we calculate:

$$\begin{aligned} (V \circ T)(v_j) &= V(T(v_j)) = V\left(\sum_{k=1}^m a_{kj} w_k\right) \\ &= \sum_{k=1}^m a_{kj} V(w_k) = \sum_{k=1}^m a_{kj} \left(\sum_{i=1}^p c_{ik} z_i\right) \\ &= \sum_{k=1}^m \sum_{i=1}^p a_{kj} c_{ik} z_i = \sum_{i=1}^p \sum_{k=1}^m c_{ik} a_{kj} z_i \\ &= \sum_{i=1}^p \left(\sum_{k=1}^m c_{ik} a_{kj} \right) z_i. \end{aligned}$$

So, the ij -entry of the matrix of $V \circ T$ is $\sum_{k=1}^m c_{ik} a_{kj}$, which is precisely the ij -entry of the matrix product CA . We thereby see that CA is the matrix of $V \circ T$ relative to the ordered bases X and Z . So, *matrix multiplication corresponds to composition of \mathbb{Z} -linear maps* (as long as the two maps use the same “middle” ordered basis Y).

15.10 Elementary Operations and Change of Basis

As in the previous section, let $T : G \rightarrow H$ be a \mathbb{Z} -linear map between two finitely generated free commutative groups G and H . Once we fix an input ordered basis X for G and an output ordered basis Y for H , we obtain a unique integer-valued matrix A representing T . It is crucial to realize that this matrix depends on the ordered bases X and Y as well as the map T .

This raises the possibility of changing the matrix of T by replacing X and Y by other ordered bases. Our eventual goal is to select ordered bases for G and H cleverly, so that the matrix of T will assume an especially simple form. Before pursuing this objective, we need to understand precisely how modifications of the ordered bases X and Y affect the matrix A .

Recall from §15.4 that there are three elementary operations on ordered bases: (B1) interchanges the positions of two basis elements; (B2) multiplies a basis element by -1 ; and (B3) adds an integer multiple of one basis element to a different basis element. We shall soon see that these operations on ordered bases correspond to analogous elementary operations

on the rows and columns of A . In particular, we consider the following elementary row operations on an integer-valued matrix (analogous to the row operations used to solve linear equations via Gaussian elimination): (R1) interchanges two rows of the matrix; (R2) multiplies some row of the matrix by -1 ; (R3) adds an integer multiple of one row to a different row. There are similar elementary operations (C1), (C2), and (C3) that act on the columns of the matrix. The key question we wish to explore is how performing the operations (B1), (B2), (B3) on Y or X causes an associated row or column operation on the matrix representing T .

The rules are most readily understood with the aid of a concrete example. Consider once again the matrix

$$A = \begin{bmatrix} 7 & 2 & -1 & 5 \\ 0 & 4 & 0 & 2 \\ 3 & 3 & 1 & 0 \end{bmatrix},$$

which is the matrix of a \mathbb{Z} -linear map $T : \mathbb{Z}^4 \rightarrow \mathbb{Z}^3$ relative to the standard ordered bases $X = (v_1, v_2, v_3, v_4)$ and $Y = (w_1, w_2, w_3)$ of \mathbb{Z}^4 and \mathbb{Z}^3 . Let us first study the effect of applying one of the basis operations (B1), (B2), and (B3) to the input basis X .

In the case of (B1), let us find the matrix of T relative to the ordered bases (v_1, v_4, v_3, v_2) and (w_1, w_2, w_3) . Here we have switched the second and fourth basis elements in X . Since $T(v_1) = 7w_1 + 0w_2 + 3w_3$, the first column of this matrix has entries 7, 0, 3. Since $T(v_4) = 5w_1 + 2w_2 + 0w_3$, the second column of this matrix has entries 5, 2, 0. Continuing in this way, we thereby obtain the matrix

$$A_1 = \begin{bmatrix} 7 & 5 & -1 & 2 \\ 0 & 2 & 0 & 4 \\ 3 & 0 & 1 & 3 \end{bmatrix}.$$

Note that A_1 is obtained from A by interchanging columns 2 and 4. One sees that this will hold in general: *if we modify the input basis by switching the vectors in positions i and j , then the new matrix is obtained by interchanging columns i and j .*

In the case of (B2), let us find the matrix of T relative to the ordered bases $(v_1, v_2, -v_3, v_4)$ and (w_1, w_2, w_3) . Here we have multiplied the third basis element in X by -1 . Evidently, columns 1, 2, and 4 of the matrix are unchanged by this modification. To find the new column 3, compute

$$T(-v_3) = -T(v_3) = w_1 + 0w_2 - w_3.$$

So the new matrix is

$$A_2 = \begin{bmatrix} 7 & 2 & 1 & 5 \\ 0 & 4 & 0 & 2 \\ 3 & 3 & -1 & 0 \end{bmatrix},$$

which was obtained from A by multiplying column 3 by -1 . This remark holds in general: *if we modify the input basis by negating the i 'th vector, then the new matrix is obtained by negating the i 'th column.*

In the case of (B3), let us find the matrix of T relative to the ordered bases $(v_1, v_2, v_3 + 2v_4, v_4)$ and (w_1, w_2, w_3) . As before, columns 1, 2, and 4 of the matrix are the same as in the original matrix A . To find the new column 3, compute

$$T(v_3 + 2v_4) = T(v_3) + 2T(v_4) = (-w_1 + w_3) + 2(5w_1 + 2w_2) = 9w_1 + 4w_2 + w_3.$$

So the new matrix is

$$A_3 = \begin{bmatrix} 7 & 2 & 9 & 5 \\ 0 & 4 & 4 & 2 \\ 3 & 3 & 1 & 0 \end{bmatrix},$$

which was obtained from A by adding two times column 4 to column 3. One can verify that this holds in general: *if we modify the i 'th input basis vector by adding c times the j 'th basis vector to it, then the new matrix is obtained by adding c times column j to column i .*

We have seen that performing elementary operations on the input basis causes elementary column operations on the matrix; we are about to see that performing elementary operations on the output basis causes elementary row operations on the matrix.

In the case of (B1), let us find the matrix of T relative to the ordered bases (v_1, v_2, v_3, v_4) and (w_1, w_3, w_2) . Taking into account the new ordering of the output basis, we have

$$T(v_1) = 7w_1 + 0w_2 + 3w_3 = 7w_1 + 3w_3 + 0w_2.$$

So the first column of the new matrix has entries 7, 3, 0 (in this order) instead of 7, 0, 3. The other columns are affected similarly. So the new matrix is

$$A_4 = \begin{bmatrix} 7 & 2 & -1 & 5 \\ 3 & 3 & 1 & 0 \\ 0 & 4 & 0 & 2 \end{bmatrix},$$

which is obtained from A by interchanging rows 2 and 3. This holds in general: *if we modify the output basis by switching the elements in positions i and j , then the associated matrix is found by interchanging rows i and j .*

In the case of (B2), let us find the matrix of T relative to the ordered bases (v_1, v_2, v_3, v_4) and $(-w_1, w_2, w_3)$. First,

$$T(v_1) = 7w_1 + 0w_2 + 3w_3 = (-7)(-w_1) + 0w_2 + 3w_3,$$

so the first column of the new matrix has entries $-7, 0, 3$. Second,

$$T(v_2) = 2w_1 + 4w_2 + 3w_3 = (-2)(-w_1) + 4w_2 + 3w_3,$$

so the second column of the new matrix has entries $-2, 4, 3$. Continuing similarly, we obtain the new matrix

$$A_5 = \begin{bmatrix} -7 & -2 & 1 & -5 \\ 0 & 4 & 0 & 2 \\ 3 & 3 & 1 & 0 \end{bmatrix},$$

which is obtained from A by multiplying the first row by -1 . This holds in general: *if we modify the output basis by negating the i 'th element, then the associated matrix is found by negating the i 'th row.*

In the case of (B3), let us find the matrix of T relative to the ordered bases (v_1, v_2, v_3, v_4) and $(w_1 + 2w_3, w_2, w_3)$. Compute

$$\begin{aligned} T(v_1) &= 7w_1 + 0w_2 + 3w_3 = 7(w_1 + 2w_3) + 0w_2 + (3 - 7 \cdot 2)w_3; \\ T(v_2) &= 2w_1 + 4w_2 + 3w_3 = 2(w_1 + 2w_3) + 4w_2 + (3 - 2 \cdot 2)w_3; \\ T(v_3) &= -1w_1 + 0w_2 + 1w_3 = -1(w_1 + 2w_3) + 0w_2 + (1 - (-1) \cdot 2)w_3; \\ T(v_4) &= 5w_1 + 2w_2 + 0w_3 = 5(w_1 + 2w_3) + 2w_2 + (0 - 5 \cdot 2)w_3. \end{aligned}$$

This leads to the new matrix

$$A_6 = \begin{bmatrix} 7 & 2 & -1 & 5 \\ 0 & 4 & 0 & 2 \\ -11 & -1 & 3 & -10 \end{bmatrix},$$

which is obtained from A by *adding -2 times row 1 to row 3*. This result, which may not

have been the one the reader was expecting, generalizes as follows: *if we modify the output basis by replacing w_i by $w_i + cw_j$, then the associated matrix is found by adding $-c$ times row i to row j .* We sketch the proof of the general case, taking $i = 1$ and $j = 2$ solely for notational convenience. Relative to the original input basis (v_1, \dots, v_n) and output basis (w_1, \dots, w_m) , we have

$$T(v_k) = a_{1k}w_1 + a_{2k}w_2 + \sum_{i=3}^m a_{ik}w_i \quad (1 \leq k \leq n).$$

We now replace w_1 by $w_1 + cw_2$. To maintain equality, we must subtract the term $ca_{1k}w_2$. Regrouping terms, we get

$$T(v_k) = a_{1k}(w_1 + cw_2) + (a_{2k} - ca_{1k})w_2 + \sum_{i=3}^m a_{ik}w_i.$$

Thus, for all k , the k 'th column of the new matrix has entries $a_{1k}, a_{2k} - ca_{1k}, a_{3k}$, etc. So, as stated in the rule, the new matrix is indeed found by adding $-c$ times row 1 to row 2.

15.11 Reduction Theorem for Integer Matrices

Now that we understand the connection between \mathbb{Z} -linear maps and integer matrices, we can concentrate our attention on the simplification of integer matrices. In this section, we will prove the following *reduction theorem for integer matrices* and then deduce implications for \mathbb{Z} -linear maps between finitely generated free commutative groups.

Reduction Theorem: Let A be an $m \times n$ matrix with integer entries. We can perform a finite sequence of elementary row and column operations on A to reduce A to a matrix of the form

$$\begin{bmatrix} a_1 & 0 & 0 & \dots & 0 \\ 0 & a_2 & 0 & \dots & 0 \\ 0 & 0 & a_3 & \dots & 0 \\ \vdots & & & \ddots & \end{bmatrix}, \quad (15.4)$$

in which there are $r \geq 0$ positive integers a_1, \dots, a_r on the main diagonal such that a_i divides a_{i+1} for all $i < r$, and all other entries in the matrix are zero. It can be shown (Exercise 87) that the matrix (15.4) satisfying the stated conditions is uniquely determined by A ; this matrix is sometimes called the *Smith normal form* of A .

We give an inductive proof of the theorem, which can be translated into an explicit recursive algorithm for reducing a given input matrix A to the required form. The base cases of the induction occur when $m = 0$ or $n = 0$ or when every entry of A is zero. In these cases, A already has the required form, so there is nothing to do.

For the induction step, assume that the reduction theorem is known to hold when the total number of rows and columns of A is less than $m + n$. Our strategy in this step is to transform A to a matrix of the form

$$\begin{bmatrix} a_1 & 0 & \dots & 0 \\ 0 & & & \\ \vdots & & A' & \\ 0 & & & \end{bmatrix}, \quad (15.5)$$

where $a_1 > 0$ and A' is an $(m - 1) \times (n - 1)$ integer-valued matrix all of whose entries are divisible by a_1 . Assuming that this has been done, we can use the induction hypothesis to continue to reduce A' to a matrix with some positive entries a_2, \dots, a_r on its diagonal, such that a_i divides a_{i+1} for $2 \leq i < r$, and with zeroes elsewhere. The operations used to reduce A' will not affect the zeroes in row 1 and column 1 of the overall matrix. Furthermore, one can check that if an integer a_1 divides every entry of a matrix, and if we perform an elementary row or column operation on that matrix, then a_1 still divides every entry of the new matrix (Exercise 52). So, as we continue to reduce A' , a_1 will always divide every entry of all the matrices obtained along the way. In particular, at the end, a_1 will divide a_2 , and we will have reduced A to a matrix of the required form.

To summarize, we need only find a way of reducing A to a matrix of the form (15.5). We assume that A is not a zero matrix, since that situation was handled in the base cases. There are two possibilities to consider. Case 1: There exists an entry $d = a_{ij}$ in A such that d divides all the entries of A . In this case, switch rows 1 and i , and then switch columns 1 and j , to bring d into the 1, 1-position. Since d divides every other entry in its column, we can subtract appropriate integer multiples of row 1 from the other rows to produce zeroes below d in column 1. Similarly, we can use column operations to produce zeroes in row 1 to the right of d . By the remark in the last paragraph, d continues to divide all the entries in the matrix as we perform these various operations on rows and columns. Finally, we can multiply row 1 by -1 to make d positive, if necessary. We now have a matrix of the form (15.5), which completes the proof in this case.

Case 2: There does not exist an entry of A that divides all the entries of A . When this case occurs, we adopt the following strategy. Let $m(A)$ be the smallest of the integers $|a_{ij}|$ as a_{ij} runs over the *nonzero* entries of A . We now aim to reduce A to a new matrix A_2 such that $m(A) > m(A_2)$. We will then repeat the whole reduction algorithm on A_2 . If A_2 satisfies case 1, then we can finish reducing as above. On the other hand, if A_2 satisfies case 2, we will reduce A_2 to a new matrix A_3 such that $m(A_2) > m(A_3)$. Continuing in this way, either case 1 will eventually occur (in which case the reduction succeeds), or case 2 occurs indefinitely. But in the latter situation, we have an infinite strictly decreasing sequence of positive integers

$$m(A) > m(A_2) > m(A_3) > \dots,$$

which violates the least natural number axiom for \mathbb{N} . So the reduction of A will always terminate after a finite number of steps.

We must still explain how to obtain the matrix A_2 such that $m(A) > m(A_2)$. For this, let $e = a_{ij}$ be a nonzero entry of A with minimal absolute value (so that $m(A) = |e|$). By definition of case 2, there exist entries in the matrix that are not divisible by e . We proceed to consider various subcases.

Case 2a: There exists $k \neq j$ such that e does not divide a_{ik} . In other words, there is an entry in e 's row that is not divisible by e . Dividing a_{ik} by e , we then have $a_{ik} = qe + r$ where $0 < r < |e|$. Subtracting q times column j from column k will produce a matrix A_2 with an r in the i, k -position. Evidently, $m(A_2) \leq r < |e|$, so we have achieved our goal in this case.

Case 2b: There exists $k \neq i$ such that e does not divide a_{kj} . In other words, there is an entry in e 's column that is not divisible by e . Dividing a_{kj} by e , we then have $a_{kj} = qe + r$ where $0 < r < |e|$. Subtracting q times row i from row k will produce a matrix A_2 with an r in the k, j -position. As before, $m(A_2) \leq r < |e|$, so we have achieved our goal in this case.

Case 2c: e divides every entry in its row and column, but for some $k \neq i, t \neq j$, e does not divide a_{kt} . So, for some $u, v \in \mathbb{Z}$, the matrix A looks like this (where we only show the

four relevant entries in rows i, k and columns j, t):

$$\left[\begin{array}{cccc} & \cdots & & \\ e & \cdots & ue & \\ \vdots & \vdots & \vdots & \vdots \\ ve & \cdots & a_{kt} & \\ & \cdots & & \end{array} \right].$$

Now, add $(1 - v)$ times row i to row k , obtaining:

$$\left[\begin{array}{cccc} & \cdots & & \\ e & \cdots & ue & \\ \vdots & \vdots & \vdots & \vdots \\ e & \cdots & a_{kt} + (1 - v)ue & \\ & \cdots & & \end{array} \right].$$

The new k, t -entry is not divisible by e , for otherwise e would divide a_{kt} . So we can proceed as in case 2a, subtracting an appropriate multiple of column j from column t to get a new k, t -entry $r < |e|$. The new matrix A_2 satisfies $m(A_2) \leq r < |e|$. The case analysis is finally complete, and the reduction theorem has now been proved. An example of the reduction process for a specific integer-valued matrix appears in §15.14 below.

15.12 Structure of \mathbb{Z} -Linear Maps between Free Groups

With the reduction theorem in hand, we can now uncover the underlying structure of \mathbb{Z} -linear maps and finitely generated commutative groups. First, suppose $T : G \rightarrow H$ is a \mathbb{Z} -linear map between two finitely generated free commutative groups. Start with any ordered bases X and Y for G and H , and let A be the matrix of T relative to X and Y . Use row and column operations to reduce A to the form given in the reduction theorem, modifying the bases X and Y as one proceeds. At the end, we will have a new ordered basis $X' = (x_1, \dots, x_n)$ for G and a new ordered basis $Y' = (y_1, \dots, y_m)$ for H such that the matrix of T relative to X' and Y' looks like

$$\left[\begin{array}{ccccc} a_1 & 0 & 0 & \dots & 0 \\ 0 & a_2 & 0 & \dots & 0 \\ 0 & 0 & a_3 & \dots & 0 \\ & & & \ddots & \end{array} \right],$$

where a_1, \dots, a_r are positive integers such that a_i divides a_{i+1} for all $i < r$. Looking at the columns of this matrix, we can give a simple description of the action of T in terms of the x_i 's and y_j 's. More specifically, inspection of the matrix shows that $T(x_i) = a_i y_i$ for $1 \leq i \leq r$ and $T(x_i) = 0$ for $r < i \leq n$. We call the matrix of T displayed above the *Smith normal form* for the \mathbb{Z} -linear map T ; it is uniquely determined by T (cf. Exercise 87).

Compare these results to the corresponding fact about linear transformations between vector spaces over fields. In that setting, we are allowed to multiply rows and columns by arbitrary nonzero scalars. The net effect of this extra ability is that we can make all the a_i 's in the final matrix become 1's (Exercise 55). The reduction process is also easier because we can use any nonzero scalar to create zeroes in all the other entries in its row and column

— no integer division is needed. So, if $T : V \rightarrow W$ is a linear map between vector spaces, there exist ordered bases $X = (x_1, \dots, x_n)$ for V and $Y = (y_1, \dots, y_m)$ for W such that $T(x_i) = y_i$ for $1 \leq i \leq r$, and $T(x_i) = 0$ for $r < i \leq n$. The number r is called the *rank* of the linear map T . One can also reach this result without reducing any matrices: start with a basis (x_{r+1}, \dots, x_n) for the null space of T , and extend it to a basis of (x_1, \dots, x_n) of V . One can then prove that the images of the first r x_i 's under T are independent, so these can be extended to a basis of the target space Y . We used this argument to prove the rank-nullity theorem in §1.8.

15.13 Structure of Finitely Generated Commutative Groups

Returning to commutative groups, we are now ready to prove the existence part of the fundamental structure theorem for finitely generated commutative groups. Suppose H is a commutative group generated by m elements. We have seen (§15.7) that H is isomorphic to a quotient group \mathbb{Z}^m/P , where P is some subgroup of the free commutative group \mathbb{Z}^m . We have also seen (§15.8) that the subgroup P must also be *free*, with a basis of $n \leq m$ elements.

Consider the inclusion map $I : P \rightarrow \mathbb{Z}^m$, given by $I(x) = x$ for $x \in P$. I is certainly \mathbb{Z} -linear, and it is a map between two finitely generated free commutative groups. So our structural result for such maps can be applied. We see, therefore, that there is a basis (x_1, \dots, x_n) for P , a basis (y_1, \dots, y_m) of \mathbb{Z}^m , an integer $r \geq 0$, and positive integers a_1, \dots, a_r with a_i dividing a_{i+1} for $i < r$, such that $I(x_i) = a_i y_i$ for $i \leq r$, and $I(x_i) = 0$ for $i > r$. Since I is an inclusion map, this says that $x_i = a_i y_i$ for $i \leq r$, and $x_i = 0$ for $i > r$. But basis elements are never zero (by \mathbb{Z} -independence), so we deduce that $r = n \leq m$. To summarize, $H \cong \mathbb{Z}^m/P$, where \mathbb{Z}^m has some ordered basis $(y_1, \dots, y_n, \dots, y_m)$ and P has an ordered basis $(a_1 y_1, \dots, a_n y_n)$.

We will now apply an isomorphism of \mathbb{Z}^m with itself that will force the subgroup P to assume the special form given in (15.2). Consider the function $f : \{y_1, \dots, y_m\} \rightarrow \mathbb{Z}^m$ such that $f(y_i) = e_i$ (the standard basis vector) for $1 \leq i \leq m$. By the UMP for free commutative groups, f extends to a \mathbb{Z} -linear map $T : \mathbb{Z}^m \rightarrow \mathbb{Z}^m$. As argued at the end of §15.6, T is an isomorphism. The isomorphism T maps P to a new subgroup P_1 of \mathbb{Z}^m with ordered basis $(a_1 e_1, \dots, a_n e_n)$. One checks that T induces an isomorphism from the quotient group \mathbb{Z}^m/P to the quotient group \mathbb{Z}^m/P_1 . Now, we can write P_1 as the Cartesian product

$$a_1 \mathbb{Z} \times a_2 \mathbb{Z} \times \cdots \times a_n \mathbb{Z} \times 0 \mathbb{Z} \times \cdots \times 0 \mathbb{Z},$$

where there are $m - n$ factors equal to $\{0\}$. Applying the fundamental homomorphism theorem as discussed below (15.2), we conclude that

$$H \cong \mathbb{Z}^m/P \cong \mathbb{Z}^m/P_1 \cong \mathbb{Z}_{a_1} \oplus \cdots \oplus \mathbb{Z}_{a_n} \oplus \mathbb{Z}^{m-n}.$$

Thus, *every finitely generated commutative group is isomorphic to a direct sum of finitely many cyclic groups, where the sizes of the finite cyclic summands (if any) successively divide each other*. Note that if some of the a_i 's are equal to 1, we can omit these factors without harm, since \mathbb{Z}_1 is the one-element group.

We now derive another version of this result, obtained by “splitting apart” the cyclic groups \mathbb{Z}_{a_i} based on the prime factorizations of the a_i . First we need a group-theoretic lemma. Suppose $a > 1$ is an integer with prime factorization $a = p_1^{e_1} p_2^{e_2} \cdots p_s^{e_s}$, where each $e_i \geq 1$ and the p_i 's are distinct primes. We claim that $\mathbb{Z}_a \cong \mathbb{Z}_{p_1^{e_1}} \oplus \mathbb{Z}_{p_2^{e_2}} \oplus \cdots \oplus \mathbb{Z}_{p_s^{e_s}}$, or

equivalently,

$$\mathbb{Z}/a\mathbb{Z} \cong (\mathbb{Z}/p_1^{e_1}\mathbb{Z}) \oplus (\mathbb{Z}/p_2^{e_2}\mathbb{Z}) \oplus \cdots \oplus (\mathbb{Z}/p_s^{e_s}\mathbb{Z}). \quad (15.6)$$

To prove this, call the product group on the right side K , and consider the map sending the integer 1 to the s -tuple of cosets $(1 + p_1^{e_1}\mathbb{Z}, \dots, 1 + p_s^{e_s}\mathbb{Z}) \in K$. Since \mathbb{Z} is free with basis $\{1\}$, the UMP furnishes a \mathbb{Z} -linear extension $T : \mathbb{Z} \rightarrow K$ such that

$$T(n) = nT(1) = (n + p_1^{e_1}\mathbb{Z}, \dots, n + p_s^{e_s}\mathbb{Z}) \quad (n \in \mathbb{Z}).$$

What is the kernel of T ? By the formula just written, n lies in the kernel iff all cosets $n + p_i^{e_i}\mathbb{Z}$ are zero iff $n \in p_i^{e_i}\mathbb{Z}$ for all i iff $p_i^{e_i}$ divides n for all i iff $\text{lcm}(p_i^{e_i} : 1 \leq i \leq s)$ divides n iff $a = p_1^{e_1} \cdots p_s^{e_s}$ divides n (since the p_i 's are distinct primes). In other words, the kernel of T is $a\mathbb{Z}$. By the fundamental homomorphism theorem, the quotient group $\mathbb{Z}/a\mathbb{Z} \cong \mathbb{Z}_a$ is isomorphic to the image of T in K . But the size of the product group K is $p_1^{e_1}p_2^{e_2} \cdots p_s^{e_s} = a$, which is the same size as $\mathbb{Z}/a\mathbb{Z}$. It follows that the isomorphic copy of $\mathbb{Z}/a\mathbb{Z}$ in K must be all of K , completing the proof of the claim.

Apply this result to each of the integers a_i in the preceding decomposition of H . We conclude that *every finitely generated commutative group is isomorphic to a direct sum of finitely many cyclic groups, each of which is either infinite or has size equal to a prime power.*

This concludes the proof of the *existence part* of the fundamental theorem of finitely generated commutative groups. We must still address the question of the *uniqueness* of the two decompositions just obtained. We will consider this issue shortly, but first we digress to give a concrete example illustrating the reduction algorithm and the ideas in the proofs just given.

15.14 Example of the Reduction Algorithm

Let P be the subgroup of \mathbb{Z}^4 generated by the vectors $v_1 = (10, 0, -8, 4)$, $v_2 = (12, 6, -6, -6)$, and $v_3 = (20, 48, 14, -82)$. Let us use the ideas in the last few sections to determine the structure of the quotient group $H = \mathbb{Z}^4/P$.

It is not hard to verify that v_1, v_2, v_3 are \mathbb{Z} -independent. (This will also follow from the calculations below; see Exercise 58.) So we can consider the matrix of the inclusion map $I : P \rightarrow \mathbb{Z}^4$ relative to the ordered basis $X = (v_1, v_2, v_3)$ of P and the standard ordered basis $Y = (e_1, e_2, e_3, e_4)$ of \mathbb{Z}^4 . The j 'th column of this matrix gives the coordinates of v_j relative to the standard ordered basis, so the matrix is

$$A = \begin{bmatrix} 10 & 12 & 20 \\ 0 & 6 & 48 \\ -8 & -6 & 14 \\ 4 & -6 & -82 \end{bmatrix}.$$

We proceed to reduce this matrix. Inspection reveals that no entry of A divides every other entry, so our first goal (following the proof of case 2 of the reduction theorem) is to reduce the magnitude of the smallest nonzero entry of A . This entry is 4, which fails to divide the entry 10 in its column. As prescribed by case 2b of the reduction proof, we replace row 1 by row 1 minus two times row 4, obtaining

$$A_1 = \begin{bmatrix} 2 & 24 & 184 \\ 0 & 6 & 48 \\ -8 & -6 & 14 \\ 4 & -6 & -82 \end{bmatrix}.$$

This is the matrix of I relative to the bases (v_1, v_2, v_3) and $(e_1, e_2, e_3, e_4 + 2e_1)$.

The new 1, 1-entry, namely 2, does divide every entry of the matrix. So we are in case 1 of the reduction proof. First, we use two column operations to produce zero entries in the rest of row 1:

$$A_2 = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 6 & 48 \\ -8 & 90 & 750 \\ 4 & -54 & -450 \end{bmatrix}.$$

This is the matrix of I relative to the bases $(v_1, v_2 - 12v_1, v_3 - 92v_1)$ and $(e_1, e_2, e_3, e_4 + 2e_1)$. Second, we use two row operations to produce zero entries in the rest of column 1:

$$A_3 = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 6 & 48 \\ 0 & 90 & 750 \\ 0 & -54 & -450 \end{bmatrix}.$$

This is the matrix of I relative to the bases $(v_1, v_2 - 12v_1, v_3 - 92v_1)$ and $(e_1 - 4e_3 + 2(e_4 + 2e_1), e_2, e_3, e_4 + 2e_1)$.

Now, we proceed to reduce the 3×2 submatrix in the lower-right corner. The upper-left entry of this submatrix (namely 6) already divides every other entry of this submatrix. Adding -8 times column 2 to column 3 gives

$$A_4 = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 90 & 30 \\ 0 & -54 & -18 \end{bmatrix};$$

the new bases are $(v_1, v_2 - 12v_1, v_3 - 92v_1 - 8(v_2 - 12v_1))$ and $(5e_1 - 4e_3 + 2e_4, e_2, e_3, e_4 + 2e_1)$. Next, two row operations produce

$$A_5 = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 30 \\ 0 & 0 & -18 \end{bmatrix};$$

the new bases are $(v_1, v_2 - 12v_1, 4v_1 - 8v_2 + v_3)$ and

$$(5e_1 - 4e_3 + 2e_4, e_2 + 15e_3 - 9(e_4 + 2e_1), e_3, e_4 + 2e_1).$$

Finally, we must reduce the 2×1 submatrix with entries 30 and -18 . We are in case 2b again; adding 2 times row 4 to row 3 gives

$$A_6 = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & -6 \\ 0 & 0 & -18 \end{bmatrix};$$

the new bases are $(v_1, v_2 - 12v_1, 4v_1 - 8v_2 + v_3)$ and

$$(5e_1 - 4e_3 + 2e_4, -18e_1 + e_2 + 15e_3 - 9e_4, e_3, e_4 + 2e_1 - 2e_3).$$

Next, multiply row 3 by -1 , obtaining the matrix

$$A_7 = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 6 \\ 0 & 0 & -18 \end{bmatrix}$$

and bases $(v_1, v_2 - 12v_1, 4v_1 - 8v_2 + v_3)$ and

$$(5e_1 - 4e_3 + 2e_4, -18e_1 + e_2 + 15e_3 - 9e_4, -e_3, e_4 + 2e_1 - 2e_3).$$

Finally, add 3 times row 3 to row 4 to get the reduced matrix

$$A_8 = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 6 \\ 0 & 0 & 0 \end{bmatrix}$$

and bases $(v_1, v_2 - 12v_1, 4v_1 - 8v_2 + v_3)$ and

$$(5e_1 - 4e_3 + 2e_4, -18e_1 + e_2 + 15e_3 - 9e_4, -e_3 - 3(e_4 + 2e_1 - 2e_3), e_4 + 2e_1 - 2e_3).$$

To check our work, note that the final ordered input basis for P is (x_1, x_2, x_3) , where:

$$\begin{aligned} x_1 &= v_1 = (10, 0, -8, 4); \\ x_2 &= v_2 - 12v_1 = (-108, 6, 90, -54); \\ x_3 &= 4v_1 - 8v_2 + v_3 = (-36, 0, 30, -18); \end{aligned}$$

and the final ordered output basis for \mathbb{Z}^4 is (y_1, y_2, y_3, y_4) , where:

$$\begin{aligned} y_1 &= 5e_1 - 4e_3 + 2e_4 = (5, 0, -4, 2); \\ y_2 &= -18e_1 + e_2 + 15e_3 - 9e_4 = (-18, 1, 15, -9); \\ y_3 &= -6e_1 + 5e_3 - 3e_4 = (-6, 0, 5, -3); \\ y_4 &= 2e_1 - 2e_3 + e_4 = (2, 0, -2, 1). \end{aligned}$$

As predicted by the proof, we have $x_1 = 2y_1$, $x_2 = 6y_2$ and $x_3 = 6y_3$. Therefore,

$$H \cong \frac{\mathbb{Z}^4}{P} \cong \frac{\mathbb{Z} \times \mathbb{Z} \times \mathbb{Z} \times \mathbb{Z}}{2\mathbb{Z} \times 6\mathbb{Z} \times 6\mathbb{Z} \times 0\mathbb{Z}} \cong \mathbb{Z}_2 \oplus \mathbb{Z}_6 \oplus \mathbb{Z}_6 \oplus \mathbb{Z}.$$

Using the prime factorization $6 = 2 \cdot 3$, we also have

$$H \cong \mathbb{Z}_2 \oplus \mathbb{Z}_2 \oplus \mathbb{Z}_2 \oplus \mathbb{Z}_3 \oplus \mathbb{Z}_3 \oplus \mathbb{Z}.$$

15.15 Some Special Subgroups

To finish our classification of finitely generated commutative groups, we want to prove the following two uniqueness theorems. First, if

$$\mathbb{Z}^b \oplus \mathbb{Z}_{a_1} \oplus \cdots \oplus \mathbb{Z}_{a_r} \cong G \cong \mathbb{Z}^d \oplus \mathbb{Z}_{c_1} \oplus \cdots \oplus \mathbb{Z}_{c_t} \quad (15.7)$$

where $b, d, r, t \geq 0$, every a_i and every c_j is > 1 , a_i divides a_{i+1} for all $i < r$, and c_j divides c_{j+1} for all $j < t$, then $b = d$ and $r = t$ and $a_i = c_i$ for $1 \leq i \leq r$. We call b the *Betti number* of G , and we call a_1, \dots, a_r the *invariant factors* of G .

Second, if

$$\mathbb{Z}^b \oplus \mathbb{Z}_{a_1} \oplus \cdots \oplus \mathbb{Z}_{a_r} \cong G \cong \mathbb{Z}^d \oplus \mathbb{Z}_{c_1} \oplus \cdots \oplus \mathbb{Z}_{c_t} \quad (15.8)$$

where $b, d, r, t \geq 0$, $a_1 \geq a_2 \geq \cdots \geq a_r > 1$, $c_1 \geq c_2 \geq \cdots \geq c_t > 1$, and every a_i and every

c_j is a prime power, then $b = d$ and $r = t$ and $a_i = c_i$ for $1 \leq i \leq r$. Again, b is called the Betti number of G ; the prime powers a_1, \dots, a_r are called *elementary divisors* of G .

We will prove these results in stages, first considering the cases where $r = t = 0$ (which means G is free) and where $b = d = 0$ (which means G is finite). To aid our proofs, we must first introduce some special subgroups of G that are invariant under group isomorphisms. Let G be an arbitrary commutative group, and let n be any fixed integer. We have a map $M_n : G \rightarrow G$ given by $M_n(g) = ng$ for $g \in G$. Because G is commutative, M_n is in fact a group homomorphism:

$$M_n(g + h) = n(g + h) = ng + nh = M_n(g) + M_n(h) \quad (g, h \in G).$$

(This result need not hold for non-commutative groups.) Define $nG = \{ng : g \in G\}$ and $G[n] = \{g \in G : ng = 0\}$. The sets nG and $G[n]$ are *subgroups* of G , since nG is the image of the homomorphism M_n , and $G[n]$ is the kernel of M_n .

We assert that these subgroups are preserved by group isomorphisms. More precisely, suppose $f : G \rightarrow H$ is a group isomorphism. We claim the restriction of f to $G[n] \subseteq G$ is a bijection of $G[n]$ onto $H[n]$. For, suppose $x \in G[n]$. Then $nx = 0$, so $f(nx) = 0$, so $nf(x) = 0$, so $f(x) \in H[n]$. So the restriction of f to the domain $G[n]$ does map into the codomain $H[n]$. Applying the same argument to the inverse isomorphism f^{-1} , we see that f^{-1} maps $H[n]$ into $G[n]$. Hence, $G[n] \cong H[n]$ via the restriction of f to this domain and codomain. The same sort of argument proves that f restricts to a group isomorphism $nG \cong nH$; here, the key point is that for $x \in nG$, we have $x = ng$ for some $g \in G$, hence $f(x) = f(ng) = nf(g)$ where $f(g) \in H$, hence $f(x) \in nH$. So f maps nG into nH , and similarly f^{-1} maps nH into nG . It now follows that f induces isomorphisms of quotient groups $G/G[n] \cong H/H[n]$ and $G/nG \cong H/nH$. The first of these isomorphisms follows, for example, by applying the fundamental homomorphism theorem to the homomorphism from G to $H/H[n]$ sending $g \in G$ to $f(g) + H[n]$, which is a surjective homomorphism with kernel $G[n]$.

Another special subgroup of G is the set $\text{tor}(G) = \bigcup_{n \geq 1} G[n]$, which consists of all elements of G of finite order: $g \in \text{tor}(G)$ iff $ng = 0$ for some $n \in \mathbb{N}^+$. We call $\text{tor}(G)$ the *torsion subgroup* of G . To see that this is a subgroup, first note that $0_G \in \text{tor}(G)$ since $1(0_G) = 0$. Given $g \in \text{tor}(G)$, we know $ng = 0$ for some $n \in \mathbb{N}^+$. Then $n(-g) = -(ng) = -0 = 0$, so $-g \in \text{tor}(G)$. To check closure under addition, suppose $g, h \in \text{tor}(G)$, so that $ng = 0 = mh$ for some $n, m > 0$. Because G is commutative,

$$nm(g + h) = nm(g) + nm(h) = m(ng) + n(mh) = m0 + n0 = 0.$$

So $g + h \in \text{tor}(G)$, and $\text{tor}(G)$ is indeed a subgroup. (Again, this result can fail for infinite non-commutative groups.) As with the previous subgroups, the torsion subgroup is an isomorphism invariant: if $f : G \rightarrow H$ is an isomorphism, one checks that f restricts to an isomorphism $\text{tor}(G) \cong \text{tor}(H)$, and so f induces an isomorphism of quotient groups $G/\text{tor}(G) \cong H/\text{tor}(H)$.

15.16 Uniqueness Proof: Free Case

To see how the subgroups introduced in the preceding section can be relevant, let us prove the uniqueness theorem for finitely generated *free* groups. Suppose G is a free commutative group with an n -element basis X and an m -element basis Y ; we will prove that $m = n$. This shows that the dimension of a finitely generated free commutative group is well-defined. The

assumptions on G imply that $G \cong \mathbb{Z}^n$ and $G \cong \mathbb{Z}^m$ (§15.6). Combining these isomorphisms gives an isomorphism $f : \mathbb{Z}^n \rightarrow \mathbb{Z}^m$. We will show that the existence of f forces $n = m$.

As we saw above, the isomorphism f induces an isomorphism $f' : \mathbb{Z}^n / 2\mathbb{Z}^n \rightarrow \mathbb{Z}^m / 2\mathbb{Z}^m$. Now, $2\mathbb{Z}^n = \{2v : v \in \mathbb{Z}^n\} = \{2(a_1, \dots, a_n) : a_i \in \mathbb{Z}\} = \{(2a_1, \dots, 2a_n) : a_i \in \mathbb{Z}\} = 2\mathbb{Z} \oplus 2\mathbb{Z} \oplus \dots \oplus 2\mathbb{Z}$. This is a subgroup of \mathbb{Z}^n of the form (15.2). So, as shown below (15.2),

$$\frac{\mathbb{Z}^n}{2\mathbb{Z}^n} = \frac{\mathbb{Z} \oplus \dots \oplus \mathbb{Z}}{2\mathbb{Z} \oplus \dots \oplus 2\mathbb{Z}} \cong \frac{\mathbb{Z}}{2\mathbb{Z}} \oplus \dots \oplus \frac{\mathbb{Z}}{2\mathbb{Z}} \cong \mathbb{Z}_2^n.$$

Similarly, $\mathbb{Z}^m / 2\mathbb{Z}^m \cong \mathbb{Z}_2^m$. So we obtain an isomorphism $\mathbb{Z}_2^n \cong \mathbb{Z}_2^m$. Now, the product group \mathbb{Z}_2^n has 2^n elements, while \mathbb{Z}_2^m has 2^m elements. Since isomorphic groups have the same size, we deduce $2^n = 2^m$, which in turn implies $n = m$.

We remark that an analogous result holds in linear algebra: any two bases of a vector space have the same cardinality. In particular, if the real vector spaces \mathbb{R}^n and \mathbb{R}^m are isomorphic, then $m = n$. However, one cannot prove this result using the above method; see Chapter 16 for a proof. Most remarkably, the *commutative groups* \mathbb{R}^n and \mathbb{R}^m are isomorphic for all positive integers m and n (see Exercise 88). Intuition suggests that the *topological spaces* \mathbb{R}^n and \mathbb{R}^m should be homeomorphic iff $m = n$. This is true, but it is quite difficult to prove when $n, m > 1$. One needs the tools of algebraic topology to establish that \mathbb{R}^n and \mathbb{R}^m are not homeomorphic when $m \neq n$. See the texts of Munkres [43] or Rotman [50] for details. We remark that algebraic topology makes heavy use of the classification theorems for commutative groups proved in this chapter.

15.17 Uniqueness Proof: Prime Power Case

For the next step in the uniqueness proof, we will prove the following result about commutative groups whose size is a prime power. *Suppose p is prime, and there are integers $a_1 \geq a_2 \geq \dots \geq a_r > 0$ and $c_1 \geq c_2 \geq \dots \geq c_t > 0$ such that*

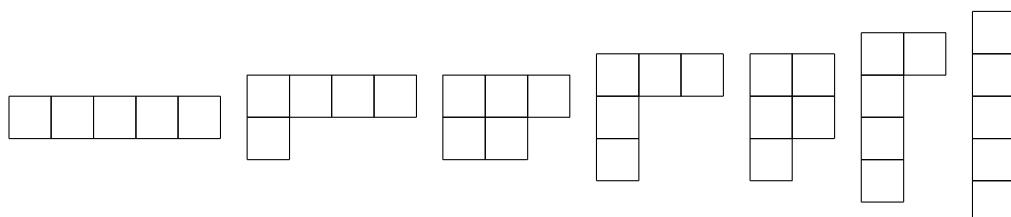
$$\mathbb{Z}_{p^{a_1}} \oplus \mathbb{Z}_{p^{a_2}} \oplus \dots \oplus \mathbb{Z}_{p^{a_r}} \cong \mathbb{Z}_{p^{c_1}} \oplus \mathbb{Z}_{p^{c_2}} \oplus \dots \oplus \mathbb{Z}_{p^{c_t}}.$$

Then $r = t$ and $a_i = c_i$ for $1 \leq i \leq r$.

The proof will be greatly clarified by introducing the notion of an integer partition (cf. Chapter 8). A *partition* of an integer n is a weakly decreasing sequence $(a_1 \geq a_2 \geq \dots \geq a_r)$ of positive integers that sum to n . For example, the seven partitions of $n = 5$ are:

$$(5), \quad (4, 1), \quad (3, 2), \quad (3, 1, 1), \quad (2, 2, 1), \quad (2, 1, 1, 1), \quad (1, 1, 1, 1, 1).$$

We can conveniently visualize a partition by drawing a collection of n squares such that there are a_i squares in row i . This picture is called the *diagram* of the partition. For example, the diagrams of the seven partitions of 5 are shown here:



Each partition encodes a possible commutative group of size p^n (where p is any fixed prime). For example, the seven partitions above correspond to the commutative groups listed here:

$$\begin{aligned} \mathbb{Z}_{p^5}, \quad & \mathbb{Z}_{p^4} \oplus \mathbb{Z}_p, \quad \mathbb{Z}_{p^3} \oplus \mathbb{Z}_{p^2}, \quad \mathbb{Z}_{p^3} \oplus \mathbb{Z}_p \oplus \mathbb{Z}_p, \\ \mathbb{Z}_{p^2} \oplus \mathbb{Z}_{p^2} \oplus \mathbb{Z}_p, \quad & \mathbb{Z}_{p^2} \oplus \mathbb{Z}_p \oplus \mathbb{Z}_p \oplus \mathbb{Z}_p, \quad \mathbb{Z}_p \oplus \mathbb{Z}_p \oplus \mathbb{Z}_p \oplus \mathbb{Z}_p \oplus \mathbb{Z}_p. \end{aligned}$$

The existence part of the fundamental structure theorem says that every commutative group of size p^5 is isomorphic to one of these seven groups. The uniqueness part of the theorem (to be proved momentarily) says that no two of these seven groups are isomorphic.

Before beginning the proof, let us see how to use partition diagrams to gain algebraic information about the associated commutative group. For definiteness, let us consider the commutative group

$$G = \mathbb{Z}_{7^4} \oplus \mathbb{Z}_{7^4} \oplus \mathbb{Z}_{7^2} \oplus \mathbb{Z}_{7^2} \oplus \mathbb{Z}_{7^2} \oplus \mathbb{Z}_7,$$

which corresponds to $p = 7$ and the partition $\mu = (4, 4, 2, 2, 2, 1)$. A typical element of G is a 6-tuple

$$x = (n_1, n_2, n_3, n_4, n_5, n_6)$$

where $0 \leq n_1 < 7^4$, $0 \leq n_2 < 7^4$, $0 \leq n_3 < 7^2$, and so on. Suppose we write n_1 as a 4-digit number in base 7:

$$n_1 = d_3 7^3 + d_2 7^2 + d_1 7^1 + d_0 7^0 \quad (0 \leq d_i < 7).$$

Writing the other n_i 's similarly, we see that we can represent the group element x by filling the squares of the partition diagram of μ with arbitrary digits in the range $\{0, 1, 2, 3, 4, 5, 6\}$. For example, the element $x = (3600, 0250, 55, 41, 30, 0)$ (where all entries of x are written in base 7) is represented by the following filled diagram:

3	6	0	0
0	2	5	0
5			
4			
3			
0			
0			

To multiply an integer written in base 7 by 7, we need only shift all the digits one place left and append a zero. In the example above, this multiplication produces

$$(36000, 02500, 550, 410, 300, 00).$$

However, to obtain $7x$ in the group G , we must now reduce the first two entries modulo 7^4 , the next three entries modulo 7^2 , and the last entry modulo 7. The effect of all these reductions is to make the leading digit of each entry disappear, leaving us with

$$7x = (6000, 2500, 50, 10, 00, 0).$$

The associated filled diagram is

6	0	0	0
2	5	0	0
5			
1			
0			
0			
0			

This diagram arises from the diagram for x by shifting the entries of each row one step left, erasing the entries that fall off the left end and bringing in zeroes at the right end. Similarly, we can compute 7^2x from the original filled diagram for x by shifting all the digits two places to the left, filling in zeroes on the right side:

$$7^2x = \begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 5 & 0 & 0 & 0 \\ \hline 0 & 0 & & \\ \hline 0 & & & \\ \hline \end{array} .$$

Inspection of the diagram shows that $7^3x = 0$, since the only nonzero digits in the diagram for x occur in the first three columns. More generally, for any $y \in G$, $7^i y = 0$ iff all the nonzero digits in the filled diagram for y occur in the first i columns of μ . This observation allows us to determine the size of the various subgroups $G[7^i] = \{y \in G : 7^i y = 0\}$. For example, how large is $G[7^2]$? Consider the following schematic diagram:

$$\begin{array}{|c|c|c|c|} \hline \star & \star & 0 & 0 \\ \hline \star & \star & 0 & 0 \\ \hline \star & \star & & \\ \hline \star & & & \\ \hline \end{array} .$$

By the previous observation, we obtain a typical element of $G[7^2]$ by choosing an arbitrary base-7 digit for each of the starred positions (possibly zero), and filling in the remaining squares of the diagram with zeroes. There are seven choices for each star, leading to the conclusion that

$$|G[7^2]| = 7^{11}.$$

These observations generalize to arbitrary diagrams. Suppose p is prime, n is a positive integer, $\mu = (a_1 \geq a_2 \geq \dots \geq a_r)$ is a partition of n , and $G = \mathbb{Z}_{p^{a_1}} \oplus \dots \oplus \mathbb{Z}_{p^{a_r}}$. Let $a'_1 \geq a'_2 \geq \dots \geq a'_s$ denote the number of squares in each column of the diagram of μ , reading from left to right. For convenience, set $a_k = 0$ for $k > r$ and $a'_k = 0$ for $k > s$. As in the example above, for any $j \in \mathbb{N}$, we obtain a typical element of $G[p^j]$ by filling the squares in the first j columns with arbitrary base- p digits, and filling the remaining squares with zeroes. By the product rule from combinatorics, we see that

$$|G[p^j]| = p^{a'_1 + a'_2 + \dots + a'_j} \quad (j \in \mathbb{N}).$$

We finally have the tools needed to attack the uniqueness proof for commutative groups of prime power size. Fix a prime p , and assume

$$G = \mathbb{Z}_{p^{a_1}} \oplus \dots \oplus \mathbb{Z}_{p^{a_r}} \cong H = \mathbb{Z}_{p^{c_1}} \oplus \dots \oplus \mathbb{Z}_{p^{c_t}},$$

where $a_1 \geq \dots \geq a_r > 0$ and $c_1 \geq \dots \geq c_t > 0$. We need to show that the two partitions $(a_i : i \geq 1)$ and $(c_i : i \geq 1)$ are the same. Let $(a'_j : j \geq 1)$ and $(c'_j : j \geq 1)$ be the column lengths of the associated partition diagrams. To show that the row lengths a_i and c_i are

the same for all i , it suffices to show that the column lengths a'_j and c'_j are the same for all j . Equivalently, it suffices to show that $a'_1 + \cdots + a'_j = c'_1 + \cdots + c'_j$ for all $j \geq 1$. To show this, note on the one hand that

$$|G[p^j]| = p^{a'_1 + \cdots + a'_j},$$

while on the other hand

$$|H[p^j]| = p^{c'_1 + \cdots + c'_j}.$$

Now $G \cong H$ implies that $G[p^j] \cong H[p^j]$, so that

$$p^{a'_1 + \cdots + a'_j} = p^{c'_1 + \cdots + c'_j}.$$

Since $p > 1$, this finally leads to $a'_1 + \cdots + a'_j = c'_1 + \cdots + c'_j$ for all j , completing the proof.

15.18 Uniqueness of Elementary Divisors

The next step is to prove the uniqueness theorems for finite commutative groups. Suppose

$$\mathbb{Z}_{a_1} \oplus \cdots \oplus \mathbb{Z}_{a_r} \cong G \cong \mathbb{Z}_{c_1} \oplus \cdots \oplus \mathbb{Z}_{c_t}, \quad (15.9)$$

where the a_i 's and c_j 's are prime powers arranged in decreasing order. We must prove $r = t$ and $a_i = c_i$ for $1 \leq i \leq r$.

The idea here is to separate out the a_j 's (resp. c_j 's) that are powers of a given prime p , so that we can apply the uniqueness result from the previous section. To see how this can be done, suppose p and q are distinct primes and $d, e \geq 1$. For $d \leq e$, every element x of the cyclic group \mathbb{Z}_{p^d} satisfies $p^e x = p^{e-d}(p^d x) = 0$, so $\mathbb{Z}_{p^d}[p^e] = \mathbb{Z}_{p^d}$. On the other hand, for arbitrary d , $\mathbb{Z}_{q^d}[p^e] = \{0\}$. For if $y \in \mathbb{Z}_{q^d}$ satisfies $p^e y = 0$, then q^d must divide the product $p^e y$ computed in \mathbb{Z} , so q^d divides y (since p and q are distinct primes). As $0 \leq y < q^d$, it follows that $y = 0$. Finally, it is routine to check that for any integer b ,

$$(G_1 \oplus \cdots \oplus G_k)[b] = (G_1[b]) \oplus \cdots \oplus (G_k[b]).$$

Now consider the situation in (15.9). Let p be any fixed prime, and let e be the highest power of p dividing $|G|$. Compute $G[p^e]$ in two ways, using the two isomorphic versions of G as direct sums of cyclic groups appearing in (15.9). On one hand, by the remarks in the previous paragraph, we will obtain a direct sum of the form

$$H_1 \oplus \cdots \oplus H_r,$$

where $H_i = \mathbb{Z}_{a_i}$ if a_i is a power of p , and $H_i = \{0\}$ if a_i is a power of another prime. Deleting zero factors and rearranging, we see that this direct sum is isomorphic to the direct sum involving precisely those a_i 's that are powers of p , arranged in decreasing order. On the other hand, applying the same reasoning to the other representation of G , we see that $G[p^e]$ is also isomorphic to the direct sum involving those c_j 's that are powers of p , arranged in decreasing order. Since $G[p^e]$ is an isomorphism invariant, the two direct sums just mentioned are isomorphic. By the results in the previous section, the a_i 's and c_j 's that are powers of the given prime p must be the same (counting multiplicities). Since p was arbitrary and every a_i and c_j is a power of some prime, we see that all the a_i 's and c_j 's must match.

15.19 Uniqueness of Invariant Factors

We can use the last result to deduce the other uniqueness theorem for finite commutative groups. Suppose now that

$$\mathbb{Z}_{a_1} \oplus \cdots \oplus \mathbb{Z}_{a_r} \cong G \cong \mathbb{Z}_{c_1} \oplus \cdots \oplus \mathbb{Z}_{c_t} \quad (15.10)$$

where the a_i 's and c_j 's are integers > 1 such that a_i divides a_{i+1} and c_j divides c_{j+1} for all $i < r$ and all $j < t$. Again, we wish to prove that $r = t$ and $a_i = c_i$ for all i . Recall (§15.13) that we can “split” and “merge” cyclic groups based on their prime factorizations, using the isomorphism

$$\mathbb{Z}_{p_1^{e_1} \cdots p_s^{e_s}} \cong \mathbb{Z}_{p_1^{e_1}} \oplus \cdots \oplus \mathbb{Z}_{p_s^{e_s}} \quad (p_k \text{'s being distinct primes}).$$

Let us use this fact to split each of the cyclic groups of size a_i and c_j in (15.10) into direct sums of cyclic groups of prime power size. Let $|G|$ have prime factorization $p_1^{e_1} \cdots p_s^{e_s}$; write $a_i = \prod_k p_k^{f_{ki}}$ and $c_j = \prod_k p_k^{g_{kj}}$ for some integers $f_{ki}, g_{kj} \geq 0$. Then

$$\bigoplus_{k=1}^s \bigoplus_{i=1}^r \mathbb{Z}_{p_k^{f_{ki}}} \cong G \cong \bigoplus_{k=1}^s \bigoplus_{j=1}^t \mathbb{Z}_{p_k^{g_{kj}}}.$$

By the previously proved uniqueness result, for each fixed k between 1 and s , the list of numbers $p_k^{f_{ki}}$ with $f_{ki} > 0$ (counted with multiplicities) equals the list of numbers $p_k^{g_{kj}}$ with $g_{kj} > 0$ (counted with multiplicities). So it suffices to show that the a_i 's (resp. c_j 's) are uniquely determined by the list of $p_k^{f_{ki}}$'s (resp. $p_k^{g_{kj}}$'s) together with the divisibility conditions linking successive a_i 's (resp. c_j 's).

We describe an algorithm for reconstructing the a_i 's from the list of all prime powers $p_k^{f_{ki}}$ that exceed 1. Construct an $s \times r$ matrix whose k 'th row contains the powers $p_k^{f_{ki}}$ arranged in increasing order, padded on the left with 1's so that there are exactly r entries in the row. Note that s is known, being the number of distinct prime factors of $|G|$. Furthermore, letting n_k be the number of positive powers of p_k that occur in the given list, r can be calculated as the maximum of n_1, \dots, n_s . This follows since $a_1 > 1$, which means that at least one prime p_k divides all r of the a_i 's. Now, every prime power in the matrix arose by splitting the prime factorization of one of the a_i 's into prime powers. Taking the divisibility relations among the a_i 's into account, it follows that a_r must be the product of the *largest* possible power of each prime. So, a_r is the product of the prime powers in the r 'th (rightmost) column of the matrix. Proceeding inductively from right to left, it follows similarly that a_{r-1} must be the product of the prime powers in column $r-1$ of the matrix, and so on. Thus we can recover a_r, a_{r-1}, \dots, a_1 uniquely from the given matrix of prime powers.

The following example of the reconstruction algorithm should clarify the preceding argument. Suppose we are given the list of prime powers:

$$[2, 2^4, 2^4, 2^4, 2^7, 3^2, 3^2, 3^2, 7, 7, 7^5, 7^8].$$

We see that $s = 3$ and $r = 5$; the matrix of prime powers is

$$\left[\begin{array}{ccccc} 2 & 2^4 & 2^4 & 2^4 & 2^7 \\ 1 & 1 & 3^2 & 3^2 & 3^2 \\ 1 & 7 & 7 & 7^5 & 7^8 \end{array} \right].$$

Multiplying the entries in each column, we get

$$a_1 = 2, \quad a_2 = 2^4 \cdot 7, \quad a_3 = 2^4 \cdot 3^2 \cdot 7, \quad a_4 = 2^4 \cdot 3^2 \cdot 7^5, \quad a_5 = 2^7 \cdot 3^2 \cdot 7^8.$$

Evidently, the a_i 's do successively divide each other, and “splitting” the a_i 's does produce the given list of prime powers. Furthermore, the reader can check for this example that the choice of a_5 (and then a_4 , etc.) is indeed forced by these conditions.

15.20 Uniqueness Proof: General Case

To finish proving the uniqueness assertion of the fundamental theorem in full generality, we use torsion subgroups to separate the “torsion part” and the “free part” of a finitely generated commutative group. More precisely, consider a group

$$G = \mathbb{Z}^b \times H,$$

where $b \geq 0$ and H is a (possibly empty) direct sum of finitely many finite cyclic groups. One sees readily that

$$\text{tor}(G) = \{0\} \times H \cong H,$$

from which it follows that

$$G/\text{tor}(G) = \frac{\mathbb{Z}^b \times H}{\{0\} \times H} \cong \frac{\mathbb{Z}^b}{\{0\}} \times \frac{H}{H} \cong \mathbb{Z}^b.$$

Let us begin the uniqueness proof. Suppose $G_1 = \mathbb{Z}^b \times H$ is isomorphic to $G_2 = \mathbb{Z}^d \times K$, where $H = \mathbb{Z}_{a_1} \times \cdots \times \mathbb{Z}_{a_r}$, $K = \mathbb{Z}_{c_1} \times \cdots \times \mathbb{Z}_{c_t}$, $b, d, r, t \geq 0$, every a_i and c_j is an integer > 1 , and either (i) a_i divides a_{i+1} and c_i divides c_{i+1} for all applicable i , or (ii) the a_i 's and c_i 's are prime powers arranged in decreasing order. We must show $b = d$ and $r = t$ and $a_i = c_i$ for all i .

Since $G_1 \cong G_2$, we know $\text{tor}(G_1) \cong \text{tor}(G_2)$ and $G_1/\text{tor}(G_1) \cong G_2/\text{tor}(G_2)$ (§15.15). In light of the preceding remarks, this means that $H \cong K$ and $\mathbb{Z}^b \cong \mathbb{Z}^d$. The isomorphism between H and K guarantees that $r = t$ and $a_i = c_i$ for all i (as shown in §15.18 and §15.19). The isomorphism between \mathbb{Z}^b and \mathbb{Z}^d guarantees that $b = d$ (by §15.16). The proof of the fundamental structure theorem for finitely generated commutative groups is finally complete.

15.21 Summary

We now review the facts about commutative groups established in this chapter.

1. *Definitions.* A *commutative group* is a set G closed under a commutative, associative binary operation (written additively) that has an identity element and additive inverses. A map $T : G \rightarrow H$ between commutative groups is \mathbb{Z} -*linear* iff $T(x + y) = T(x) + T(y)$ for all $x, y \in G$ (which automatically implies $T(nx) = nT(x)$ for all $x \in G$ and $n \in \mathbb{Z}$). A \mathbb{Z} -*linear combination* of elements $v_1, \dots, v_k \in G$ is any element of the form $n_1v_1 + \cdots + n_kv_k$ where the n_i 's are integers. G is *generated* by v_1, \dots, v_k iff every $g \in G$ is a \mathbb{Z} -linear combination of the v_i 's. The v_i 's are \mathbb{Z} -*independent* iff $n_1v_1 + \cdots + n_kv_k = 0$ ($n_i \in \mathbb{Z}$) implies all n_i 's are zero. The v_i 's are a \mathbb{Z} -*basis* of G iff they generate G and are \mathbb{Z} -independent iff for each $g \in G$ there are *unique* integers n_i with $g = \sum_i n_i v_i$. A commutative

group is *finitely generated* iff it has a finite generating set; the group is *free* iff it has a basis; the group is *k-dimensional* iff it has a basis with k elements.

2. *Properties of Generating Sets.* Two \mathbb{Z} -linear maps $S, T : G \rightarrow H$ that agree on a generating set of G must be equal. A \mathbb{Z} -linear map $T : G \rightarrow H$ is surjective iff its image contains a generating set for H .
3. *Universal Mapping Property for Free Commutative Groups.* Suppose X is a \mathbb{Z} -basis of a commutative group G . Given any commutative group H and any function $f : X \rightarrow H$, there exists a unique \mathbb{Z} -linear extension $T_f : G \rightarrow H$ such that $T_f(v) = f(v)$ for all $v \in X$.
4. *Consequences of the Universal Mapping Property.* Every k -dimensional free commutative group G is isomorphic to \mathbb{Z}^k . Different isomorphisms can be obtained by choosing an ordered basis (v_1, \dots, v_k) for G and sending $n_1 v_1 + \dots + n_k v_k \in G$ to the k -tuple of coordinates $(n_1, \dots, n_k) \in \mathbb{Z}^k$. Every commutative group with a k -element generating set is isomorphic to a quotient group of \mathbb{Z}^k .
5. *Matrix Representation of \mathbb{Z} -Linear Maps.* Given a \mathbb{Z} -linear map $T : G \rightarrow H$, an ordered basis $X = (v_1, \dots, v_n)$ of G , and an ordered basis $Y = (w_1, \dots, w_m)$ of H , the *matrix of T relative to the input basis X and output basis Y* is the unique $m \times n$ integer-valued matrix $A = [a_{ij}]$ such that

$$T(v_j) = \sum_{i=1}^m a_{ij} w_i \quad (1 \leq j \leq n).$$

Matrix addition corresponds to pointwise addition of \mathbb{Z} -linear maps, while matrix multiplication corresponds to composition of linear maps.

6. *Elementary Operations on Bases, Rows, and Columns.* Given an ordered \mathbb{Z} -basis of a free commutative group G , we can create new ordered bases of G by the following operations: switch two basis vectors; negate a basis vector; add an integer multiple of one basis vector to another basis vector. Similar operations can be performed on the rows and columns of integer matrices. Column operations on the matrix of a \mathbb{Z} -linear map T correspond to changes in the input basis, while row operations correspond to changes in the output basis.
7. *Reduction of Integer Matrices.* Using row and column operations, we can reduce any $m \times n$ integer-valued matrix A to a new matrix B such that the nonzero entries of B (if any) occupy the first r positions on the main diagonal of B , and these nonzero entries are positive integers each of which divides the next one. B is uniquely determined by A and is called the *Smith normal form* of A .
8. *Canonical Form for \mathbb{Z} -Linear Maps.* Suppose $T : G \rightarrow H$ is a \mathbb{Z} -linear map between two finite-dimensional free commutative groups. There exist an ordered basis $X = (x_1, \dots, x_n)$ for G , an ordered basis $Y = (y_1, \dots, y_m)$ for H , an integer $r \geq 0$, and positive integers a_1, \dots, a_r such that: a_i divides a_{i+1} for all $i < r$, $T(x_i) = a_i y_i$ for $1 \leq i \leq r$, and $T(x_i) = 0$ for $r < i \leq n$.
9. *Subgroups of Finitely Generated Free Commutative Groups.* Any subgroup P of a k -dimensional free commutative group G is also free with dimension at most k . By choosing appropriate bases for P and G , one can find an isomorphism $G \cong \mathbb{Z}^k$ such that P maps under this isomorphism to a subgroup of the form

$$a_1 \mathbb{Z} \oplus \dots \oplus a_m \mathbb{Z} \oplus \{0\} \oplus \dots \oplus \{0\},$$

where each a_i divides the next one. It follows that G/P is isomorphic to a direct sum of cyclic groups.

10. *Invariant Subgroups.* For any commutative group G and integer n , the subsets $G[n] = \{g \in G : ng = 0\}$, $nG = \{ng : g \in G\}$, and $\text{tor}(G) = \bigcup_{n \geq 1} G[n]$ are subgroups of G . Any isomorphism $f : G \rightarrow H$ restricts to give isomorphisms $G[n] \cong H[n]$, $nG \cong nH$, and $\text{tor}(G) \cong \text{tor}(H)$, and f induces isomorphisms of quotient groups $G/G[n] \cong H/H[n]$, $G/nG \cong H/nH$, and $G/\text{tor}(G) \cong H/\text{tor}(H)$.

11. *Fundamental Theorem of Finitely Generated Commutative Groups (Version 1).* Every finitely generated commutative group G is isomorphic to a direct sum of cyclic groups

$$\mathbb{Z}^b \oplus \mathbb{Z}_{a_1} \oplus \cdots \oplus \mathbb{Z}_{a_r},$$

where $b, r \geq 0$, every $a_i > 1$, and a_i divides a_{i+1} for all $i < r$. The integers b, r , and a_1, \dots, a_r (subject to the stated conditions) are uniquely determined by G .

12. *Fundamental Theorem of Finitely Generated Commutative Groups (Version 2).* Every finitely generated commutative group G is isomorphic to a direct sum of cyclic groups

$$\mathbb{Z}^b \oplus \mathbb{Z}_{a_1} \oplus \cdots \oplus \mathbb{Z}_{a_r},$$

where $b, r \geq 0$, $a_1 \geq \cdots \geq a_r > 1$, and every a_i is a prime power. The integers b, r , and a_1, \dots, a_r (subject to the stated conditions) are uniquely determined by G .

13. *Restatement of Uniqueness in the Fundamental Theorem.* Suppose

$$G = \mathbb{Z}^b \oplus \mathbb{Z}_{a_1} \oplus \cdots \oplus \mathbb{Z}_{a_r} \cong \mathbb{Z}^d \oplus \mathbb{Z}_{c_1} \oplus \cdots \oplus \mathbb{Z}_{c_t},$$

where $b, d, r, t \geq 0$, all a_i 's and c_j 's are > 1 , and either: (i) the a_i 's and c_j 's are prime powers arranged in decreasing order; or (ii) a_i divides a_{i+1} for all $i < r$ and c_j divides c_{j+1} for all $j < t$. Then $b = d$ and $r = t$ and $a_i = c_i$ for $1 \leq i \leq r$. In particular (when $r = t = 0$), this says that the dimension of a finite-dimensional free commutative group is well-defined. The a_i 's in case (i) are called the *elementary divisors* of G . The a_i 's in case (ii) are called the *invariant factors* of G .

14. *Remarks for Non-Finitely Generated Commutative Groups.* A subset S of a commutative group G generates G iff every $g \in G$ is a \mathbb{Z} -linear combination of some finite subset of S ; S is \mathbb{Z} -independent iff every finite list of distinct elements of S is \mathbb{Z} -independent; S is a \mathbb{Z} -basis of G iff S is \mathbb{Z} -independent and generates G . A commutative group G is called *free* iff it has a \mathbb{Z} -basis. The additive group \mathbb{Q} is an example of a commutative group that is not finitely generated. The universal mapping property is valid for free commutative groups. Every free commutative group G is isomorphic to a direct sum of copies of \mathbb{Z} , where the number of factors in the direct sum is the cardinality of a basis of G . Every commutative group is isomorphic to a quotient group of a free commutative group.

15.22 Exercises

Unless otherwise specified, assume $(G, +)$ and $(H, +)$ are commutative groups in these exercises.

1. (a) For fixed $n \in \mathbb{Z}$, verify that $n\mathbb{Z} = \{ni : i \in \mathbb{Z}\}$ is a subgroup of \mathbb{Z} . (b) Find all $m, n \in \mathbb{Z}$ with $m\mathbb{Z} = n\mathbb{Z}$.

2. Let H be a subgroup of G . Prove: for all $x \in G$, $\langle x \rangle \subseteq H$ iff $x \in H$.
3. Fix $n \in \mathbb{N}^+$. (a) Prove that every element of $\mathbb{Z}/n\mathbb{Z}$ equals one of the cosets $0 + n\mathbb{Z}$, $1 + n\mathbb{Z}$, \dots , $(n - 1) + n\mathbb{Z}$, and show that these cosets are all distinct. (b) Prove that $\mathbb{Z}/n\mathbb{Z}$ is isomorphic to (\mathbb{Z}_n, \oplus) as groups. (c) Prove that $\mathbb{Z}/0\mathbb{Z} \cong \mathbb{Z}$.
4. In the text, we constructed a cyclic group of size $n \in \mathbb{N}^+$ by forming the quotient group $\mathbb{Z}/n\mathbb{Z}$. (a) Find a specific subgroup of (S_n, \circ) that is isomorphic to $\mathbb{Z}/n\mathbb{Z}$. (b) Find a specific subgroup of \mathbb{C}^* (the nonzero complex numbers under multiplication) that is isomorphic to $\mathbb{Z}/n\mathbb{Z}$. (c) Find a specific subgroup of $\text{GL}_n(\mathbb{R})$ (real invertible $n \times n$ matrices under matrix multiplication) that is isomorphic to $\mathbb{Z}/n\mathbb{Z}$.
5. Let G_1, \dots, G_k be commutative groups. (a) Check that $G = G_1 \times \dots \times G_k$ is a commutative group under componentwise addition. (b) For $1 \leq i \leq k$, let $j_i : G_i \rightarrow G$ be given by $j_i(x_i) = (0, \dots, x_i, \dots, 0)$ for $x_i \in G_i$, where the x_i occurs in position i . Check that each j_i is an injective group homomorphism. (c) Which results in (a) and (b) are true for arbitrary groups G_1, \dots, G_k ?
6. Let $n \in \mathbb{N}$ and $x \in G$. The definition of nx in §15.2 can be stated more carefully in the following recursive way: $0x = 0$, and $(n+1)x = (nx) + x$. Use this recursive definition and induction to prove the following “laws of multiples” from (15.1) assuming $m, n \in \mathbb{N}$ and $x, y \in G$ (a separate argument is needed when m or n is negative): (a) $(m+n)x = (mx) + (nx)$; (b) $m(nx) = (mn)x$; (c) $1x = x$; (d) $m(x+y) = (mx) + (my)$.
7. Give a specific example of a non-commutative group $(G, +)$, $x, y \in G$, and $m \in \mathbb{N}^+$ with $m(x+y) \neq mx+my$.
8. Let $v_1, \dots, v_k \in G$, and let $H = \{c_1v_1 + \dots + c_kv_k : c_i \in \mathbb{Z}\}$. (a) Prove that the identity $0_G \in H$ and that H is closed under inverses. (b) If $(G, +)$ is not commutative, must H be closed under the group operation? Prove or give a counterexample. (c) Let K be the set of all group elements of the form $e_1w_1 + e_2w_2 + \dots + e_sw_s$, where $s \in \mathbb{N}$, each $e_i \in \{1, -1\}$, and each w_i is some v_j (different w_i 's could be equal to the same v_j). Prove that for all G (commutative or not), K is a subgroup of G , and prove $K = H$ when G is commutative.
9. Find all generators of each of these groups: (a) \mathbb{Z}_8 ; (b) \mathbb{Z}_{12} ; (c) \mathbb{Z}_{30} ; (d) $\mathbb{Z}_3 \times \mathbb{Z}_5$.
10. (a) Prove that for prime p and all nonzero a in \mathbb{Z}_p , $\mathbb{Z}_p = \langle a \rangle$. (b) For prime p and $e \in \mathbb{N}^+$, how many $a \in \mathbb{Z}_{p^e}$ satisfy $\mathbb{Z}_{p^e} = \langle a \rangle$?
11. For $n \in \mathbb{N}^+$ and $k \in \mathbb{Z}_n$, find and prove a criterion for when $\mathbb{Z}_n = \langle k \rangle$.
12. (a) Suppose G_1 and G_2 are commutative groups with $G_1 = \langle v_1, \dots, v_m \rangle$ and $G_2 = \langle w_1, \dots, w_n \rangle$. Prove $G_1 \times G_2 = \langle (v_1, 0), \dots, (v_m, 0), (0, w_1), \dots, (0, w_n) \rangle$. (b) Generalize (a) to finite direct products $G_1 \times G_2 \times \dots \times G_k$.
13. (a) Describe an explicit generating set for the product group $G = \mathbb{Z} \times \dots \times \mathbb{Z} \times \mathbb{Z}_{a_1} \times \dots \times \mathbb{Z}_{a_s}$, where there are b copies of \mathbb{Z} and every $a_i = p_i^{e_i}$ is a prime power. (b) How many generating sets S does G have in which $|S| = b+s$, and every element of S has exactly one nonzero component?
14. Let S be an infinite subset of G . (a) Prove the set H of \mathbb{Z} -linear combinations of elements of S is a subgroup of G . (b) Prove that $S = \{1/p^e : p \text{ is prime}, e \in \mathbb{N}^+\}$ generates $(\mathbb{Q}, +)$. (c) Prove or disprove: the set $T = \{1/p : p \text{ is prime}\}$ generates $(\mathbb{Q}, +)$. (d) Find a generating set for \mathbb{Q}^+ , the group of positive rational numbers under multiplication.

15. (a) Show that $\mathbb{Z}[x]$ (polynomials with integer coefficients, under addition) is a commutative group that is not finitely generated. Find an infinite generating set for this group. (b) Give an example of a commutative group $(G, +)$ that is not finitely generated and is not isomorphic to $\mathbb{Z}[x]$ or to $(\mathbb{Q}, +)$. (c) Prove or disprove: the commutative groups $\mathbb{Z}[x]$ and $(\mathbb{Q}, +)$ are isomorphic. (d) Prove or disprove: the commutative groups $\mathbb{Z}[x]$ and \mathbb{Q}^+ (positive rationals under multiplication) are isomorphic.
16. Let $B = (v_1, \dots, v_k)$ be a list of elements in G . (a) Say what it means for B to be \mathbb{Z} -dependent, by negating the definition of \mathbb{Z} -linear independence. (b) Say what it means for a subset S of G (possibly infinite) to be \mathbb{Z} -dependent. (c) Use (b) to explain why \emptyset is \mathbb{Z} -independent relative to G .
17. Which of the following commutative groups are free? Explain. (a) $3\mathbb{Z}$ (a subgroup of \mathbb{Z}); (b) $\mathbb{Z} \times \mathbb{Z}_5$; (c) $\mathbb{Z}[x]$ under addition; (d) $\{a + bi : a \in 2\mathbb{Z}, 2b \in \mathbb{Z}\}$ under complex addition; (e) \mathbb{Q} under addition.
18. Let V be a \mathbb{Q} -vector space. Prove that a list (v_1, \dots, v_k) of elements of V is \mathbb{Z} -linearly independent iff the list is \mathbb{Q} -linearly independent.
19. (a) Show that for all $a \in \mathbb{N}^+$ with a not a perfect square, $(1, \sqrt{a})$ is a \mathbb{Z} -linearly independent list in \mathbb{R} . (b) Show that $(1, \sqrt{2}, \sqrt{3}, \sqrt{6})$ is \mathbb{Z} -linearly independent.
20. Let $B = (v_1, \dots, v_k)$ be a list of elements in G . Suppose C is obtained from B by switching v_i and v_j , for some $i \neq j$. (a) Prove $G = \langle B \rangle$ iff $G = \langle C \rangle$. (b) Prove B is \mathbb{Z} -independent iff C is \mathbb{Z} -independent. (c) Prove B is an ordered basis of G iff C is an ordered basis of G .
21. (a) Repeat Exercise 20, but now assume C is obtained from B by replacing some v_i by $-v_i$. (b) Which implications in Exercise 20 are true if C is obtained from B by replacing some v_i by nv_i , where $n \notin \{-1, 0, 1\}$ is a fixed integer?
22. Check that $B = ((0, 1, -5), (0, 0, -1), (1, 2, 3))$ is an ordered basis of \mathbb{Z}^3 : (a) by proving from the definitions that B generates \mathbb{Z}^3 and is \mathbb{Z} -linearly independent; (b) by showing how to obtain B from the known ordered basis (e_1, e_2, e_3) by a sequence of elementary operations.
23. Repeat Exercise 22 for $B = ((12, -7, -2), (8, 5, 3), (1, 2, 1))$.
24. (a) Find and prove necessary and sufficient conditions on $a, b, c, d \in \mathbb{Z}$ so that $((a, b), (c, d))$ is an ordered basis of \mathbb{Z}^2 . (b) Can you generalize your answer to (a) to characterize ordered bases of \mathbb{Z}^k for all $k \geq 1$?
25. Let F be a field, and let V be an F -vector space with ordered basis $B = (v_1, \dots, v_n)$. Prove that the three elementary operations in §15.4 (as modified in the last paragraph of that section) send the basis B to another ordered basis of V .
26. Prove or disprove: for all $k \geq 1$, every ordered basis of the free commutative group \mathbb{Z}^k can be obtained from the standard ordered basis (e_1, e_2, \dots, e_k) by a finite sequence of elementary operations (B1), (B2), and (B3).
27. (a) Give a justified example of a finitely generated free commutative group G and a \mathbb{Z} -linearly independent list (v_1, \dots, v_k) in G that cannot be extended to an ordered \mathbb{Z} -basis (v_1, \dots, v_s) of G . (b) Give a justified example of a finitely generated free commutative group G and a generating set $\{v_1, \dots, v_k\}$ of G such that no subset of this generating set is a \mathbb{Z} -basis of G .
28. Let G and H be commutative groups, and let $T : G \rightarrow H$ be a group homomorphism. (a) Prove by induction on n that $T(nx) = nT(x)$ for all $x \in G$

- and all $n \in \mathbb{N}$. (b) Prove $T(nx) = nT(x)$ for all $x \in G$ and all $n \in \mathbb{Z}$. (c) Prove by induction on k that $T(\sum_{i=1}^k n_i v_i) = \sum_{i=1}^k n_i T(v_i)$ for all $k \in \mathbb{N}$, $n_i \in \mathbb{Z}$ and $v_i \in G$.
29. Find the coordinates of $(1, 2, 3)$ and $(4, -1, 1)$ relative to each ordered basis for \mathbb{Z}^3 : (a) $B_1 = (e_1, e_2, e_3)$; (b) $B_2 = ((0, 1, -5), (0, 0, -1), (1, 2, 3))$; (c) $B_3 = ((12, -7, -2), (8, 5, 3), (1, 2, 1))$.
 30. Let $T : G \rightarrow H$ be a \mathbb{Z} -linear map. (a) Prove: If $X \subseteq G$ generates G , then $T[X] = \{T(x) : x \in X\}$ generates $\text{img}(T)$. (b) Prove or disprove: If $X \subseteq G$ generates G and $Y = X \cap \ker(T)$, then Y generates $\ker(T)$.
 31. Let G be a free commutative group with ordered basis $B = (v_1, v_2, \dots, v_k)$. (a) Use the UMP for free commutative groups to construct an isomorphism $T : G \rightarrow \mathbb{Z}v_1 \oplus \mathbb{Z}v_2 \oplus \dots \oplus \mathbb{Z}v_k$. (b) Show that $\mathbb{Z}v_i \cong \mathbb{Z}$, preferably by using the UMP, and conclude that $G \cong \mathbb{Z}^k$.
 32. Let G and H be finitely generated free commutative groups of dimensions n and m , respectively. Use the UMP for free commutative groups to prove the following facts: (a) If $n \leq m$, there exists an injective \mathbb{Z} -linear map $S : G \rightarrow H$. (b) If $n \geq m$, there exists a surjective \mathbb{Z} -linear map $T : G \rightarrow H$. (c) If $n = m$, there exists a bijective \mathbb{Z} -linear map $U : G \rightarrow H$.
 33. (a) Give a justified example of a free commutative group G , a commutative group H , a finite spanning set X for G , and a function $f : X \rightarrow H$ that has no extension to a \mathbb{Z} -linear map with domain G . (b) Give a justified example of a free commutative group G , a commutative group H , a \mathbb{Z} -independent subset X of G , and a function $f : X \rightarrow H$ that has infinitely many extensions to \mathbb{Z} -linear maps with domain G .
 34. Let G be a free commutative group with infinite basis X . (a) Prove for every $g \in G$, there exist unique integers $\{n_x : x \in X\}$ with only finitely many n_x 's nonzero and $g = \sum_{x \in X} n_x x$. (b) Prove that if $S, T : G \rightarrow H$ are \mathbb{Z} -linear maps that agree on X , then $S = T$. (c) Prove for every commutative group H and every function $f : X \rightarrow H$, there exists a unique \mathbb{Z} -linear map $T_f : G \rightarrow H$ extending f . (d) Check that for any set Y , the set $\mathbb{Z}^{(Y)}$ of all functions $f : Y \rightarrow \mathbb{Z}$ such that $f(y) = 0$ for all but finitely many $y \in Y$ is a free commutative group under pointwise addition of functions, with a basis in bijective correspondence with the set Y . (e) Prove G is isomorphic to the group $\mathbb{Z}^{(X)}$. (f) Prove that every commutative group is isomorphic to a quotient group of a free commutative group.
 35. Fix $k, n \in \mathbb{N}^+$, and let $X = \{e_1, \dots, e_k\}$ where each $e_i = (0, \dots, 1, \dots, 0)$ is viewed as an element of $G = \mathbb{Z}_n^k$. Prove that G and X satisfy the following UMP: for all commutative groups H such that $ny = 0$ for all $y \in H$ and for all functions $f : X \rightarrow H$, there exists a unique \mathbb{Z} -linear map $T_f : G \rightarrow H$ extending f .
 36. Use the UMP for the free commutative group \mathbb{Z} , together with the fundamental homomorphism theorem for groups, to show that every cyclic group is isomorphic to \mathbb{Z} or to $\mathbb{Z}/n\mathbb{Z}$ for some $n \in \mathbb{N}^+$.
 37. Check that the map S defined below (15.2) is a surjective \mathbb{Z} -linear map with kernel P .
 38. Let G_1, \dots, G_k be groups with respective normal subgroups H_1, \dots, H_k . Prove
- $$(G_1 \times \dots \times G_k)/(H_1 \times \dots \times H_k) \cong (G_1/H_1) \times \dots \times (G_k/H_k).$$
39. Suppose G is a free commutative group, A and B are commutative groups,

$f : A \rightarrow B$ is a surjective group homomorphism, and $g : G \rightarrow B$ is a group homomorphism. Prove there exists a group homomorphism $h : G \rightarrow A$ with $f \circ h = g$.

40. In §15.8, we proved that every subgroup of a finitely generated free commutative group is free. (a) Trace through the construction in that proof, applied to the subgroup $H = \{(t, 2t) : t \in \mathbb{Z}\}$ of \mathbb{Z}^2 , to construct an ordered basis for this subgroup. (b) Similarly, use the proof to find an ordered basis for $H = \langle(2, 4, 15), (4, 6, 6)\rangle$, which is a subgroup of \mathbb{Z}^3 .
41. Let F be a field. Modify the proof in §15.8 to prove that for all $k \in \mathbb{N}$, every subspace of F^k has an ordered basis of size at most k . Use only the definitions and induction, avoiding any theorems whose conclusions involve the existence of a basis.
42. Let $e_1 = (1, 0)$, $e_2 = (0, 1)$, $f_1 = (2, 5)$, $f_2 = (1, 3)$. Let $T : \mathbb{Z}^2 \rightarrow \mathbb{Z}^2$ be the \mathbb{Z} -linear map given by $T((a, b)) = (4a - b, 2a + 3b)$ for $a, b \in \mathbb{Z}$. (a) Find the matrix of T relative to the input basis (e_1, e_2) and output basis (e_1, e_2) . (b) Find the matrix of T relative to the input basis (e_1, e_2) and output basis (f_1, f_2) . (c) Find the matrix of T relative to the input basis (f_1, f_2) and output basis (e_1, e_2) . (d) Find the matrix of T relative to the input basis (f_1, f_2) and output basis (f_1, f_2) .
43. Let $v_1 = (7, 2, 2)$, $v_2 = (2, -1, 0)$, and $v_3 = (3, 1, 1)$; let $I : \mathbb{Z}^3 \rightarrow \mathbb{Z}^3$ be the identity map; and let $X = (e_1, e_2, e_3)$ and $Y = (v_1, v_2, v_3)$. (a) Check that Y is an ordered \mathbb{Z} -basis of \mathbb{Z}^3 . (b) Find the matrix of I relative to the input basis X and output basis X . (c) Find the matrix of I relative to the input basis Y and output basis X . (d) Find the matrix of I relative to the input basis X and output basis Y . (e) Find the matrix of I relative to the input basis Y and output basis Y .
44. Let $T, U : \mathbb{Z}^3 \rightarrow \mathbb{Z}^3$ be the \mathbb{Z} -linear maps whose matrices (using the standard ordered basis (e_1, e_2, e_3) as both input and output basis) are

$$A = \begin{bmatrix} 2 & 0 & -1 \\ 4 & 4 & 1 \\ 0 & 3 & -2 \end{bmatrix} \text{ and } B = \begin{bmatrix} -1 & 5 & 7 \\ 2 & 0 & -2 \\ -3 & 1 & 4 \end{bmatrix},$$

respectively. (a) Find $T((a, b, c))$ and $U((a, b, c))$ for all $a, b, c \in \mathbb{Z}$. (b) Find the matrix of $T + U$ relative to the standard ordered basis. (c) Find the matrix of TU relative to the standard ordered basis. (d) Find the matrix of UT relative to the standard ordered basis.

45. In this problem, view elements of \mathbb{Z}^n and \mathbb{Z}^m as column vectors. (a) Let $T : \mathbb{Z}^n \rightarrow \mathbb{Z}^m$ be a \mathbb{Z} -linear map with matrix A relative to the standard ordered bases of \mathbb{Z}^n and \mathbb{Z}^m . Prove: for $v \in \mathbb{Z}^n$, $T(v) = Av$ (the matrix-vector product of A and v). (b) More generally, suppose $T : G \rightarrow H$ is \mathbb{Z} -linear, $X = (x_1, \dots, x_n)$ is an ordered basis of G , $Y = (y_1, \dots, y_m)$ is an ordered basis of H , and A is the matrix of T relative to these bases. Show that if $g \in G$ has coordinates $v \in \mathbb{Z}^n$ relative to X , then $Av \in \mathbb{Z}^m$ gives the coordinates of $T(g)$ relative to Y .
46. Let G and H be finitely generated free commutative groups with ordered bases X and Y , respectively. Suppose $T : G \rightarrow H$ is a \mathbb{Z} -linear map whose matrix relative to X and Y is A . (a) For any $c \in \mathbb{Z}$, show that $cT : G \rightarrow H$, defined by $(cT)(x) = cT(x)$ for $x \in G$, is \mathbb{Z} -linear. (b) Find the matrix of cT relative to X and Y .

47. *Dual Groups.* For any commutative group G , define the *dual group* of G , denoted G^* or $\text{Hom}_{\mathbb{Z}}(G, \mathbb{Z})$, to be the set of all \mathbb{Z} -linear maps $f : G \rightarrow \mathbb{Z}$. (a) Show that G^* is a commutative group under pointwise addition of functions. (b) Show that $\mathbb{Z}^* \cong \mathbb{Z}$. (c) For commutative groups G_1, \dots, G_k , show that $(G_1 \times \cdots \times G_k)^* \cong (G_1^*) \times \cdots \times (G_k^*)$. (d) Show that if G is a free commutative group of dimension k , then G^* is a free commutative group of dimension k .
48. (a) Suppose G is free with ordered basis $X = (v_1, \dots, v_n)$. Show there exists a unique *dual basis* $X^* = (v_1^*, \dots, v_n^*)$ of the dual group G^* (defined in Exercise 47) satisfying $v_i^*(v_j) = 1$ if $i = j$ and 0 if $i \neq j$. (b) Suppose H is also free with ordered basis $Y = (w_1, \dots, w_m)$. Given a \mathbb{Z} -linear map $T : G \rightarrow H$, show the map $T^* : H^* \rightarrow G^*$ given by $T^*(f) = f \circ T$ for $f \in H^*$ is \mathbb{Z} -linear. (c) How is the matrix of T^* relative to the bases Y^* and X^* related to the matrix A of T relative to the bases X and Y ?
49. Assume the setup in Exercise 44. For each input basis X and output basis Y , find the matrix of T relative to these bases. (a) $X = (e_1 + 4e_3, e_2, e_3)$, $Y = (e_1, e_2, e_3)$; (b) $X = (e_1, e_2, -e_3)$, $Y = (e_1, e_2, e_3)$; (c) $X = (e_1, e_2, e_3)$, $Y = (e_1, e_3, e_2)$; (d) $X = (e_1, e_2, e_3)$, $Y = (-e_1, -e_2, e_3)$; (e) $X = (e_2, e_1, e_3)$, $Y = (e_1, e_2, e_3)$; (f) $X = (e_1, e_2, e_3)$, $Y = (e_1, e_2 - 3e_3, e_3)$; (g) $X = (e_1, e_2 - 2e_1, e_3 + e_1)$, $Y = (e_1 + e_3, e_2 - 2e_3, e_3)$.
50. Suppose G is a free commutative group with ordered basis (v_1, v_2, v_3, v_4) , H is a free commutative group with ordered basis (w_1, w_2, w_3) , and $T : G \rightarrow H$ is a \mathbb{Z} -linear map whose matrix (relative to these bases) is

$$A = \begin{bmatrix} 2 & -4 & -3 & 0 \\ 1 & -2 & 3 & 3 \\ 4 & -3 & -2 & 4 \end{bmatrix}.$$

- (a) Compute $T(3v_2 - v_3 + 2v_4)$. (b) What is the matrix of T relative to the input basis (v_3, v_1, v_2, v_4) and output basis (w_3, w_2, w_1) ? (c) What is the matrix of T relative to the input basis $(v_1, -v_2, v_3, v_4)$ and output basis $(-w_1, w_2, -w_3)$? (d) What is the matrix of T relative to the input basis $(v_1, v_2 - v_1, v_3 + 3v_1, v_4 + v_1)$ and output basis (w_1, w_2, w_3) ? (e) What is the matrix of T relative to the input basis (v_1, v_2, v_3, v_4) and output basis $(w_1, w_1 + w_2 + w_3, 2w_2 + w_3)$?
51. Carefully prove the six italicized statements in §15.10, which indicate how elementary operations on input and output bases affect the matrix of a \mathbb{Z} -linear map.
52. Suppose $A \in M_{m,n}(\mathbb{Z})$ and an integer b divides every entry of A . Let C be obtained from A by applying a single elementary row or column operation. (a) Show that b divides every entry of C . (b) Show that the gcd of all entries of A equals the gcd of all entries of C (interpret the gcd as zero in the case of a zero matrix). (c) Show that the 1, 1-entry of the Smith normal form of A is the gcd of all entries of A . (This result is generalized in Exercise 87.)
53. Use elementary row and column operations in \mathbb{Z} to reduce each integer matrix below to its Smith normal form (15.4):
- (a) $A = \begin{bmatrix} 9 & 18 \\ 36 & 6 \end{bmatrix}$; (b) $B = \begin{bmatrix} 9 & 8 & 7 \\ 6 & 5 & 4 \\ 3 & 2 & 1 \end{bmatrix}$; (c) $C = \begin{bmatrix} 12 & 0 & 9 & 15 \\ -21 & -9 & 27 & 18 \\ 0 & 15 & 33 & -21 \end{bmatrix}$.
54. Let $T : \mathbb{Z}^2 \rightarrow \mathbb{Z}^2$ be the \mathbb{Z} -linear map from Exercise 42. Find \mathbb{Z} -bases $X = (g_1, g_2)$ and $Y = (h_1, h_2)$ of $\mathbb{Z} \times \mathbb{Z}$ such that the matrix of T with respect to the input

basis X and the output basis Y has the form $\begin{bmatrix} c & 0 \\ 0 & d \end{bmatrix}$, where $c, d \in \mathbb{N}$ and c divides d . (Reduce the matrix found in Exercise 42(a), keeping track of how each operation changes the input or output basis.)

55. Let F be a field. Prove: For any matrix $A \in M_{m,n}(F)$, we can perform finitely many elementary row and column operations on A to obtain a matrix B such that, for some r with $0 \leq r \leq \min(m, n)$, $B(i, i) = 1_F$ for $1 \leq i \leq r$ and all other entries of B are zero. (Imitate the first part of the proof of the reduction theorem for integer-valued matrices.)
56. Write a computer program that takes as input a matrix $A \in M_{m,n}(\mathbb{Z})$ and returns as output the Smith normal form (15.4) of A .
57. Prove that a \mathbb{Z} -linear map between two finitely generated free commutative groups is invertible iff the Smith normal form for the map is an identity matrix.
58. Prove that the columns of $A \in M_{m,n}(\mathbb{Z})$ are \mathbb{Z} -linearly dependent iff the Smith normal form of A has at least one column of zeroes.
59. Let M be the subgroup of \mathbb{Z}^3 generated by $v_1 = (6, 6, 9)$ and $v_2 = (12, 6, 6)$.
 - (a) Explain why (v_1, v_2) is an ordered \mathbb{Z} -basis of M .
 - (b) Define $T : M \rightarrow \mathbb{Z}^3$ by $T(x) = x$ for $x \in M$. Find the matrix of T relative to the input basis (v_1, v_2) and output basis (e_1, e_2, e_3) .
 - (c) Use matrix reduction to find a new basis (x_1, x_2) for M , a new basis (y_1, y_2, y_3) for \mathbb{Z}^3 , and positive integers d_1, d_2 (where d_1 divides d_2) such that $x_1 = d_1 y_1$ and $x_2 = d_2 y_2$.
 - (d) Find a product of cyclic groups (satisfying the conclusions of the classification theorem for finitely generated commutative groups) that is isomorphic to \mathbb{Z}^3/M .
60. (a) In the proof in §15.13, verify that the map T extending the function f sending each y_i to e_i is an isomorphism.
- (b) Verify that $P_1 = T[P]$ has ordered basis $(a_1 e_1, \dots, a_n e_n)$.
- (c) Verify that T induces a group isomorphism from \mathbb{Z}^m/P to \mathbb{Z}^m/P_1 .
61. Let m and n be relatively prime positive integers.
 - (a) Use (15.6) to prove that $\mathbb{Z}/(mn\mathbb{Z}) \cong (\mathbb{Z}/m\mathbb{Z}) \times (\mathbb{Z}/n\mathbb{Z})$.
 - (b) Imitate the proof of (15.6) to prove that $\mathbb{Z}/(mn\mathbb{Z}) \cong (\mathbb{Z}/m\mathbb{Z}) \times (\mathbb{Z}/n\mathbb{Z})$.
 - (c) Which steps in the proof in (b) fail when $\gcd(m, n) > 1$?
62. Assume the setup in Exercise 44.
 - (a) Use matrix reduction to find new bases X and Y for \mathbb{Z}^3 such that the matrix of T relative to X and Y is in Smith normal form.
 - (b) Repeat (a) for the map U .
63. A certain \mathbb{Z} -linear map $T : \mathbb{Z}^4 \rightarrow \mathbb{Z}^3$ is represented by the following matrix relative to the standard ordered bases of \mathbb{Z}^4 and \mathbb{Z}^3 :

$$A = \begin{bmatrix} 15 & 0 & -10 & 20 \\ 30 & -20 & 20 & 10 \\ 25 & -15 & 55 & 40 \end{bmatrix}.$$

- (a) Give a formula for $T((i, j, k, p))$, where $i, j, k, p \in \mathbb{Z}$.
- (b) Use matrix reduction to find an ordered basis $X = (v_1, v_2, v_3, v_4)$ of \mathbb{Z}^4 , an ordered basis $Y = (w_1, w_2, w_3)$ of \mathbb{Z}^3 , and a matrix B in Smith normal form such that B is the matrix of T relative to X and Y .
64. Use matrix reduction to determine the Betti numbers, elementary divisors, and invariant factors for each of the following quotient groups:
 - (a) $\mathbb{Z}^2/\langle(-4, 4), (-8, -4)\rangle$;
 - (b) $\mathbb{Z}^3/\langle(-255, -12, -60), (-114, -6, -27)\rangle$;
 - (c) $\mathbb{Z}^4/\langle(50, 160, 70, 210), (69, 213, 81, 282), (29, 88, 31, 117)\rangle$.

65. Give a specific example of a group (G, \star) and a positive integer n such that the map $M_n : G \rightarrow G$ given by $M_n(g) = g^n$ for $g \in G$ is not a group homomorphism, the image of M_n is not a subgroup of G , and the set $\{g \in G : M_n(g) = e_G\}$ is not a subgroup of G .
66. Suppose $f : G \rightarrow H$ is a group homomorphism. (a) Prove: for all $n \in \mathbb{N}^+$, $f[G[n]] \subseteq H[n]$. (b) Give an example where strict inclusion holds in (a). (c) If f is injective, must equality hold in (a)? Explain. (d) If f is surjective, must equality hold in (a)? Explain.
67. Suppose $f : G \rightarrow H$ is a group homomorphism. (a) Prove: for all $n \in \mathbb{N}^+$, $f[nG] \subseteq nH$. (b) Give an example where strict inclusion holds in (a). (c) If f is injective, must equality hold in (a)? Explain. (d) If f is surjective, must equality hold in (a)? Explain.
68. Assume $|G| = m$. Prove: for all $n \in \mathbb{N}^+$ with $\gcd(m, n) = 1$, $nG = G$ and $G[n] = \{0\}$.
69. Assume $f : G \rightarrow H$ is a group homomorphism. (a) Prove $f[\text{tor}(G)] \subseteq \text{tor}(H)$. (b) Give an example where strict inclusion holds in (a). (c) Prove: if f is an isomorphism, then f restricts to an isomorphism $\text{tor}(G) \cong \text{tor}(H)$. (d) Deduce from (c) that for an isomorphism f , we get an induced isomorphism $G/\text{tor}(G) \cong H/\text{tor}(H)$.
70. In §15.16, we proved that $\mathbb{Z}^n \cong \mathbb{Z}^m$ implies $n = m$. Where does this proof break down if we try to use it to show that $\mathbb{R}^n \cong \mathbb{R}^m$ (isomorphism of additive groups) implies $n = m$?
71. For $1 \leq n \leq 4$, list the partitions of n and draw their diagrams.
72. Let p be a fixed prime. Use partitions to make a complete list of all non-isomorphic commutative groups of size p^6 .
73. Let G be the group

$$\mathbb{Z}_{5^5} \times \mathbb{Z}_{5^5} \times \mathbb{Z}_{5^5} \times \mathbb{Z}_{5^3} \times \mathbb{Z}_{5^3} \times \mathbb{Z}_{5^2} \times \mathbb{Z}_{5^2}.$$

Draw pictures of partition diagrams to help answer the following questions. (a) For each $i \geq 1$, find the size of the subgroup $G[5^i]$. (b) For each $i \geq 1$, find the size of the subgroup 5^iG . (c) Find the size of $G[125] \cap 25G$ (explain).

74. Suppose p is prime and

$$G = \mathbb{Z}_{p^{a_1}} \times \mathbb{Z}_{p^{a_2}} \times \cdots \times \mathbb{Z}_{p^{a_k}}$$

where $\mu = (a_1 \geq a_2 \geq \cdots \geq a_k)$ is a partition. (a) What is the size of the subgroup pG ? (b) Describe how to use partition diagrams to find the size of p^iG for all $i \geq 1$. (c) Describe how to use partition diagrams to find the size of $p^iG \cap G[p^j]$ for all $i, j \in \mathbb{N}$.

75. Prove that two integer partitions $(a_i : i \geq 1)$ and $(c_i : i \geq 1)$ are equal iff the column sums $a'_1 + \cdots + a'_j$ and $c'_1 + \cdots + c'_j$ are equal for all $j \geq 1$.
76. Let G_1, \dots, G_k be commutative groups and $b \in \mathbb{N}^+$. (a) Prove $(G_1 \oplus \cdots \oplus G_k)[b] = (G_1[b]) \oplus \cdots \oplus (G_k[b])$. (b) Prove $b(G_1 \oplus \cdots \oplus G_k) = (bG_1) \oplus \cdots \oplus (bG_k)$. (c) Prove $\text{tor}(G_1 \oplus \cdots \oplus G_k) = \text{tor}(G_1) \oplus \cdots \oplus \text{tor}(G_k)$.
77. For each n below, make a complete list of all non-isomorphic commutative groups of size n (use decompositions that display the elementary divisors of each group): (a) $n = 400$; (b) $n = 1001$; (c) $n = 666$; (d) $n = p^2q^3$, where p and q are distinct primes.

78. For each n below, make a complete list of all non-isomorphic commutative groups of size n (use decompositions that display the invariant factors of each group):
 (a) $n = 24$; (b) $n = 300$; (c) $n = 32$; (d) $n = p^3q^3$, where p and q are distinct primes.
79. For each commutative group, find its elementary divisors. (a) \mathbb{Z}_{9900} ;
 (b) $\mathbb{Z}_{60} \times \mathbb{Z}_{100} \times \mathbb{Z}_{80}$; (c) $\mathbb{Z}_{48} \times \mathbb{Z}_{111} \times \mathbb{Z}_{99} \times \mathbb{Z}_{1001}$.
80. For each commutative group, find its invariant factors.
 (a) $\mathbb{Z}_{32} \times \mathbb{Z}_{16} \times \mathbb{Z}_4 \times \mathbb{Z}_4 \times \mathbb{Z}_9 \times \mathbb{Z}_9 \times \mathbb{Z}_3$; (b) $\mathbb{Z}_{60} \times \mathbb{Z}_{100} \times \mathbb{Z}_{80}$; (c) $\mathbb{Z}_{48} \times \mathbb{Z}_{111} \times \mathbb{Z}_{99} \times \mathbb{Z}_{1001}$.
81. Let p and q be distinct primes. How many non-isomorphic commutative groups of size p^5q^3 are there?
82. Let P be a logical property such that every cyclic group has property P ; if a group G has property P and $G \cong H$, then H has property P ; and whenever groups G and H have property P , the product group $G \times H$ has property P . Prove that all finitely generated commutative groups have property P .
83. Use the classification of finite commutative groups to characterize all positive integers n such that every commutative group of size n is cyclic. Give two proofs, one based on elementary divisors and one based on invariant factors.
84. (a) Use the classification of finite commutative groups to prove that for all n -element commutative groups G and all positive divisors d of n , G has a subgroup of size d . (b) Prove that if G is an n -element commutative group that has at most one subgroup of size d , for each positive divisor d of n , then G must be cyclic.
 (c) Can you prove (a) without using the classification results from this chapter?
85. (a) Suppose A , B , and C are finitely generated commutative groups such that $A \times C \cong B \times C$. Prove that $A \cong B$. (b) Give an example to show that the result of (a) can fail if C is not finitely generated.
86. (a) Let G be a commutative group of size $p_1^{e_1}p_2^{e_2} \cdots p_k^{e_k}$, where p_1, \dots, p_k are distinct primes and $e_1, \dots, e_k \in \mathbb{N}^+$. Prove: for $1 \leq i \leq k$, G has a unique subgroup P_i of size $p_i^{e_i}$. (b) Give an example to show that (a) can fail if G is not commutative. (c) Show that if the conclusion in (a) holds, then $G \cong P_1 \times P_2 \times \cdots \times P_k$ even if G is not commutative. (Use Exercise 9 from Chapter 1.)
87. Let $A \in M_{m,n}(\mathbb{Z})$. For $1 \leq k \leq \min(m, n)$, a *minor* of A of order k is the determinant of some $k \times k$ submatrix of A obtained by looking at the entries in k fixed rows and k fixed columns of A . Let $G_k(A)$ be the gcd of all k 'th order minors of A (use the convention that $\gcd(0, 0, \dots, 0) = 0$). (a) Show that performing one elementary row or column operation on A does not change any of the integers $G_k(A)$. (First show that if $c \in \mathbb{Z}$ divides all k 'th order minors of A , then c divides all k 'th order minors of the new matrix.) (b) Show that if B is any Smith normal form of A , then $G_k(B) = G_k(A)$. (c) Show that if B is any matrix in Smith normal form, then $G_k(B) = \prod_{i=1}^k B(i, i)$ for $1 \leq k \leq \min(m, n)$. (d) Use (b) and (c) to prove that the Smith normal form of a matrix is unique. (e) Use (b) and (c) to compute the Smith normal form of

$$A = \begin{bmatrix} 10 & 8 & 0 \\ -4 & 8 & 6 \\ 0 & 12 & -8 \end{bmatrix}.$$

88. Let B be a basis for \mathbb{R} , viewed as a \mathbb{Q} -vector space. (a) For $m \in \mathbb{N}^+$, use B to

- describe a basis B_m for the \mathbb{Q} -vector space \mathbb{R}^m . (b) Argue that $|B| = |B_m|$ for all $m \in \mathbb{N}^+$. Conclude that all of the commutative groups $(\mathbb{R}^m, +)$ (for $m = 1, 2, 3, \dots$) are isomorphic to the commutative group $(\mathbb{R}, +)$.
89. Give justified examples of each of the following: (a) an infinite commutative group that is not free; (b) a free commutative group that is not infinite; (c) a \mathbb{Z} -independent list in \mathbb{R}^2 that is not \mathbb{R} -independent; (d) a free commutative group with an infinite basis; (e) an infinite commutative group G such that for all $x \in G$, there exists $n \in \mathbb{N}$ with $nx = 0_G$; (f) a commutative group G and a proper subgroup H with $G \cong H$; (g) a commutative group G with isomorphic subgroups A and B such that G/A is not isomorphic to G/B ; (h) a group G in which $\text{tor}(G)$ is not a subgroup of G .
90. True or false? Explain each answer. (a) Every subgroup of a finitely generated free commutative group is free. (b) Every quotient group of a finitely generated free commutative group is free. (c) A product of finitely many finitely generated free commutative groups is free. (d) If (x_1, x_2, x_3) is an ordered basis for a commutative group G , $(x_3, x_2 + x_3, x_1 + x_2 + x_3)$ must also be an ordered basis for G . (e) If (x_1, \dots, x_n) is any \mathbb{Z} -independent list in a commutative group G and c is a nonzero integer, then (cx_1, x_2, \dots, x_n) must also be \mathbb{Z} -independent. (f) For all commutative groups G , $(G \sim \text{tor}(G)) \cup \{e_G\}$ is a subgroup of G . (g) For all commutative groups G and all $n \in \mathbb{N}^+$, $G/G[n] \cong nG$. (h) Every finite commutative group is isomorphic to a subgroup of \mathbb{Z}_n for some $n \in \mathbb{N}^+$. (i) Every finite commutative group is isomorphic to a quotient group of \mathbb{Z}_n for some $n \in \mathbb{N}^+$. (j) Every finite commutative group is isomorphic to a subgroup of S_n for some $n \in \mathbb{N}^+$. (k) Every finitely generated commutative group is isomorphic to a subgroup of \mathbb{Z}^n for some $n \in \mathbb{N}^+$. (l) Every finitely generated commutative group is isomorphic to a quotient group of \mathbb{Z}^n for some $n \in \mathbb{N}^+$. (m) Every finitely generated commutative group G with $\text{tor}(G) = \{0\}$ must be free. (n) Every k -element generating set for \mathbb{Z}^k must be a \mathbb{Z} -basis for \mathbb{Z}^k . (o) Every k -element \mathbb{Z} -linearly independent subset of \mathbb{Z}^k must generate \mathbb{Z}^k . (p) For all finite commutative groups G and H , if $G \times G \cong H \times H$, then $G \cong H$.

This page intentionally left blank

Axiomatic Approach to Independence, Bases, and Dimension

The goal of this chapter is to present a general axiomatic treatment of some fundamental concepts from linear algebra: linear independence, linear dependence, spanning sets, and bases. One benefit of the axiomatic approach is that it sweeps away a lot of irrelevant extra structure, isolating a few key properties that underlie the basic theorems about linear independence and bases. Even better, these axioms arise in other situations besides linear algebra. Hence, all the theorems deduced from the axioms will apply to those other situations as well. For example, we will prove the main theorems about the transcendence degree of field extensions by verifying the axioms of the general theory.

To motivate our axiomatic treatment of dependence relations and bases, consider a vector space V over a field F . For each subset S of V , we have the subspace of V spanned by S , denoted $\text{span}_F(S)$, consisting of all F -linear combinations of finitely many elements of S . Our axioms will single out certain properties of the *spanning operator* that maps each subset S to the subspace spanned by S . From these properties, we will eventually prove the key theorems that every vector space V has a basis, and any two bases of V have the same cardinality. We will return to this example in §16.9 after developing the general theory.

16.1 Axioms

Before stating the axioms, we recall some notation from set theory that will be used extensively in this chapter. Given sets S and T , we write $S \cup T = \{x : x \in S \text{ or } x \in T\}$; $S \cap T = \{x : x \in S \text{ and } x \in T\}$; and $S \sim T = \{x : x \in S \text{ and } x \notin T\}$. To reduce notational clutter, given $x \in X$ and $S \subseteq X$, we define $S + x = S \cup \{x\}$ and $S - x = S \sim \{x\}$.

Our axiomatic framework consists of a fixed set X and a *spanning operator* Sp satisfying the following axioms:

- A0. For each subset S of X , there is a uniquely determined subset $\text{Sp}(S) \subseteq X$.
- A1. For all $S \subseteq X$, $S \subseteq \text{Sp}(S)$.
- A2. For all $S, T \subseteq X$, if $S \subseteq \text{Sp}(T)$ then $\text{Sp}(S) \subseteq \text{Sp}(T)$.
- A3. For all $x, y \in X$ and $S \subseteq X$, if $x \in \text{Sp}(S + y)$ and $x \notin \text{Sp}(S)$, then $y \in \text{Sp}(S + x)$.
- A4. For all $x \in X$ and $S \subseteq X$, if $x \in \text{Sp}(S)$ then $x \in \text{Sp}(S')$ for some finite subset S' of S .

Axiom A3 is called the *exchange axiom*. The pair (X, Sp) is called an *independence structure* or a *finitary matroid*. If we replace axiom A4 by the stronger condition that X is a finite set, we obtain one definition of a *matroid* (other definitions are discussed in §16.14 below).

16.2 Definitions

We now define abstract versions of the following concepts from linear algebra: spanning, subspaces, linear dependence, linear independence, bases, and finite-dimensionality.

- D1. For $S, V \subseteq X$, S spans V (or generates V) iff $\text{Sp}(S) = V$.
- D2. For $V \subseteq X$, V is a subspace iff $V = \text{Sp}(S)$ for some $S \subseteq X$.
- D3. For $S \subseteq X$, S is dependent iff there exists $x \in S$ with $x \in \text{Sp}(S - x)$.
- D4. For $S \subseteq X$, S is independent iff S is not dependent iff for all $x \in S$, $x \notin \text{Sp}(S - x)$.
- D5. For $S \subseteq X$, S is a basis for X iff S is independent and $\text{Sp}(S) = X$.
- D6. X is finitely generated (or finite-dimensional) iff $X = \text{Sp}(S)$ for some finite set $S \subseteq X$.

Given a logical property P , we say that a subset S of X is a maximal subset satisfying P iff P is true for S , and for any set $T \subseteq X$ properly containing S , P is not true for T . If X is finite and at least one subset of X has property P , then there exists at least one maximal subset of X with property P ; it suffices to pick a subset S of maximum possible size among all subsets satisfying property P . However, in our treatment, we are not assuming X is finite. So more care is needed when trying to find maximal subsets satisfying various properties. One tool for doing so is *Zorn's lemma*, discussed in §16.6 below.

16.3 Initial Theorems

We begin by proving some basic results that follow from axioms A0, A1, and A2. Unless otherwise stated, S and T are arbitrary subsets of X .

- T1. *The empty set \emptyset is independent.*
Proof: If not, there exists $x \in \emptyset$ (satisfying certain properties, namely $x \in \text{Sp}(\emptyset - x)$).
This is impossible, since no x is a member of the empty set.
- T2. $\text{Sp}(X) = X$, so X is a subspace.
Proof: First, $\text{Sp}(X) \subseteq X$ (by A0).
Second, $X \subseteq \text{Sp}(X)$ (let $S = X$ in A1).
- T3. *If $S \subseteq T$, then $\text{Sp}(S) \subseteq \text{Sp}(T)$.* (monotonicity of spanning operator)
Proof: We have $T \subseteq \text{Sp}(T)$ (by A1).
Therefore, $S \subseteq \text{Sp}(T)$ (since $S \subseteq T$).
So $\text{Sp}(S) \subseteq \text{Sp}(T)$ (by A2).
- T4. $\text{Sp}(\text{Sp}(T)) = \text{Sp}(T)$. (*idempotence* of spanning operator)
Proof: First, $\text{Sp}(T) \subseteq \text{Sp}(\text{Sp}(T))$ (let $S = \text{Sp}(T)$ in A1).
Second, $\text{Sp}(\text{Sp}(T)) \subseteq \text{Sp}(T)$ (let $S = \text{Sp}(T)$ in A2).
- T5. *If $S \subseteq T$ and S is dependent, then T is dependent.*
Proof: Choose $x \in S$ with $x \in \text{Sp}(S - x)$ (by D3).
We have $S - x \subseteq T - x$ (since $S \subseteq T$).
Therefore, $x \in \text{Sp}(T - x)$ (by T3).
Since $x \in T$, T is dependent (by D3).

T6. If $S \subseteq T$ and T is independent, then S is independent.

Proof: Take the contrapositive of T5.

T7. If M is a basis of X , then M is a maximal independent subset of X .

Proof: M is independent and spans X (by D5).

Let N be any subset of X properly containing M . We must show N is dependent.

Fix $x_0 \in N$ with $x_0 \notin M$.

Let $S = M + x_0$. Then S is a subset of N , $x_0 \in S$, and $S - x_0 = M$.

We have $x_0 \in \text{Sp}(M)$ (M spans X).

So $x_0 \in \text{Sp}(S - x_0)$, and S is dependent (by D3).

Therefore N , which contains S , is dependent (by T5).

The converse of theorem T7 will be proved later (in T13), with the help of the exchange axiom A3.

16.4 Consequences of the Exchange Axiom

The proofs so far have only invoked axioms A0, A1, and A2. To establish the key facts regarding existence of bases and uniqueness of the cardinality of bases, we need to use the exchange axiom A3. First, we list four logically equivalent formulations of this axiom (here $x, y \in X$ and $S \subseteq X$ are arbitrary):

A3(a). The following three conditions cannot all be true at once:

$$x \notin \text{Sp}(S); x \in \text{Sp}(S + y); y \notin \text{Sp}(S + x).$$

A3(b). If $x \notin \text{Sp}(S)$ and $x \in \text{Sp}(S + y)$, then $y \in \text{Sp}(S + x)$.

A3(c). If $x \notin \text{Sp}(S)$ and $y \notin \text{Sp}(S + x)$, then $x \notin \text{Sp}(S + y)$.

A3(d). If $x \in \text{Sp}(S + y)$ and $y \notin \text{Sp}(S + x)$, then $x \in \text{Sp}(S)$.

Second, we prove an analogous collection of four equivalent statements involving the notion of dependence.

T8. For all $y \in X$ and $U \subseteq X$, the following three conditions cannot all be true at once:

$$y \notin \text{Sp}(U); V = U + y \text{ is dependent}; U \text{ is independent}.$$

Proof: Assume all three conditions are true; we will contradict axiom A3(a).

Since V is dependent but U is not, $V \neq U$, and so $y \notin U$.

By dependence of V , choose $x \in V$ with $x \in \text{Sp}(V - x)$.

If $x = y$, then $V - x = U$, so $y = x \in \text{Sp}(U)$, contradicting our assumption.

Thus, $x \neq y$, and $x \in U$. Choose $S = U - x$.

Note that $S + x = U$ and $S + y = V - x$. Therefore:

1. $x \notin \text{Sp}(S)$ (by independence of U).

2. $x \in \text{Sp}(S + y)$ (by choice of x).

3. $y \notin \text{Sp}(S + x)$ (by assumption on y).

We have contradicted axiom A3(a).

T9. For all $y \in X$ and $U \subseteq X$, if $y \notin \text{Sp}(U)$ and $U + y$ is dependent, then U is dependent.

T10. For all $y \in X$ and $U \subseteq X$, if $y \notin \text{Sp}(U)$ and U is independent, then $U + y$ is independent.

- T11. For all $y \in X$ and $U \subseteq X$, if U is independent and $U + y$ is dependent, then $y \in \text{Sp}(U)$.

Third, we prove a strong converse to theorem T7. Together with theorem T10, this converse will be the crucial lemma for proving that bases of X exist.

- T12. If T spans X and M is a maximal independent subset of T , then M is a basis of X .

Proof: First we show $T \subseteq \text{Sp}(M)$. Let z be any element of T .

If $z \in M$, then $z \in \text{Sp}(M)$ (by A1).

Otherwise, $S = M + z$ is a subset of T properly containing M (as $z \notin M$).

Therefore S is dependent (by maximality of M).

Since $M + z$ is dependent and M is independent, $z \in \text{Sp}(M)$ (by T11).

We have now shown that $T \subseteq \text{Sp}(M)$.

Hence, $X = \text{Sp}(T) \subseteq \text{Sp}(M)$ (by A2).

Since also $\text{Sp}(M) \subseteq X$ (by A0), we have $X = \text{Sp}(M)$.

As M is independent, it is a basis of X (by D5).

- T13. If M is a maximal independent subset of X , then M is a basis of X .

Proof: Let $T = X$ in T12, keeping in mind T2.

Fourth, we prove a lemma that describes an exchange property for independent sets. This lemma will be used to compare the size of any independent set S in X to any spanning set T in X , from which we will deduce the uniqueness of the cardinality of a basis of X . The idea of the lemma is that we can replace an element of S by a carefully chosen element of T and still have an independent set. (See Figure 16.1.)

- T14. Assume S, T are subsets of X with S independent and $X = \text{Sp}(T)$.

For any $s \in S \sim T$, there exists $t \in T \sim S$ such that $(S - s) + t$ is independent.

Proof: Fix the independent set S , the spanning set T , and $s \in S \sim T$.

Since S is independent, $s \notin \text{Sp}(S - s)$.

Note $X = \text{Sp}(T) \subseteq \text{Sp}((S \cup T) - s) \subseteq X$ (by T3 and A0).

So $X = \text{Sp}((S \cup T) - s)$, and hence $s \in \text{Sp}((S \cup T) - s)$.

To get a contradiction, assume that $(T \sim S) \subseteq \text{Sp}(S - s)$.

Since $(S - s) \subseteq \text{Sp}(S - s)$ (by A1), it follows that $(S \cup T) - s \subseteq \text{Sp}(S - s)$.

Then $X = \text{Sp}((S \cup T) - s) \subseteq \text{Sp}(S - s)$ (by A2).

This is impossible, since $s \in X$ and $s \notin \text{Sp}(S - s)$.

So there exists $t \in T \sim S$ with $t \notin \text{Sp}(S - s)$.

Note $S - s$ is independent (by T6).

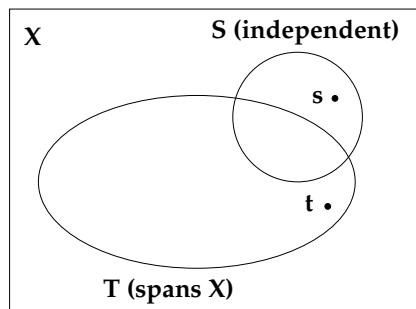
So $(S - s) + t$ is independent (by T10).

16.5 Main Theorems: Finite-Dimensional Case

For simplicity, we first prove the main results on bases in the case where X is finitely generated. The theorems in this section are not needed in the treatment of the general case, which is given in §16.7.

- T15. Assume $X = \text{Sp}(S)$ with S finite. Any finite independent set in X can be extended to a finite basis of X by adding appropriate elements of S .

Proof: Let U be any finite independent set in X .

**FIGURE 16.1**

Replacing $s \in S \sim T$ by an Appropriate $t \in T \sim S$ Preserves Independence of S .

The set $T = S \cup U$ is finite, and U is an independent subset of T containing U . Therefore, we can choose a maximal independent subset M of T containing U . The set T spans X , since its subset S spans X (T3).

Therefore, M is a finite basis of X (by T12).

Since $U \subseteq M \subseteq S \cup U$, all elements of $M \sim U$ come from S .

- T16. *Assume X is finitely generated. Any finite spanning set T of X contains a basis of X .*

So X has a finite basis.

Proof: Apply T15 to the spanning set T and the independent set \emptyset (see T1).

- T17. *Assume X is finitely generated. If S is any independent set in X and U is any spanning set for X , then $|S| \leq |U|$.*

Proof: Let S be any independent set in X , possibly infinite. Using the hypothesis of finite generation, let T be a fixed finite spanning set for X having the minimum possible cardinality among all spanning sets for X , say $|T| = n$. Since $|T| \leq |U|$, it suffices to prove that $|S| \leq |T|$.

The idea of the proof is to use T14 to show that we can replace elements of S not in T by certain distinct elements of T , one at a time, and still have an independent set. We will prove that this replacement process can always continue as long as S contains unreplaced elements. On the other hand, the process must end when the elements of T not in S have all been used up. These two facts will imply that S cannot have more elements than T .

To implement this idea, we use induction to create a sequence of independent sets $S = S_0, S_1, S_2, \dots, S_k$ with $k \leq n - |S \cap T| < \infty$, such that $|S| = |S_0| = |S_1| = |S_2| = \dots = |S_k|$ and $|S_i \cap T| = |S_0 \cap T| + i$ for $0 \leq i \leq k$. To begin, define $S_0 = S$. For the induction step, fix $i \geq 0$, and assume we have found an independent set S_i with $|S_i| = |S_0|$ and $|S_i \cap T| = |S_0 \cap T| + i$. Consider two cases. If $S_i \sim T = \emptyset$, then $S_i \subseteq T$. We stop the inductive construction at this stage, noting that $|S| = |S_i| \leq |T|$ gives the required conclusion. In the other case, $S_i \sim T$ is nonempty; say $s_i \in S_i \sim T$. Then T14 applies to S_i, T , and s_i to give $t_i \in T \sim S_i$ such that $S_{i+1} = (S_i - s_i) + t_i$ is independent. Since $s_i \in S_i \sim T$ and $t_i \in T \sim S_i$, it is evident that $|S_{i+1}| = |S_i|$ (even if S_i is infinite). We also see that deleting s_i from S_i and adding t_i increases the size of the finite set $|S_i \cap T|$ by 1, so $|S_{i+1} \cap T| = |S_i \cap T| + 1 = |S_0 \cap T| + i + 1$ as needed. The size of the intersection starts at $|S \cap T|$, increases by 1 at each step, and can never exceed

$n = |T|$. So the second case can occur at most $n - |S \cap T| < \infty$ times. Thus the first case happens after finitely many steps, completing the proof.

- T18. Assume X is finitely generated. Every independent set in X and every basis of X is finite.

Proof: This follows from T17 and the fact that a basis is an independent set.

- T19. Assume X is finitely generated. If A and B are any two bases of X , then $|A| = |B| < \infty$.

Proof: Both A and B are finite, by T18.

Since A is independent and B spans X , $|A| \leq |B|$ by T17.

Since B is independent and A spans X , $|B| \leq |A|$ by T17.

Therefore, $|A| = |B| < \infty$.

16.6 Zorn's Lemma

We want to prove versions of the theorems of the preceding section (excluding T18 and conclusions involving finiteness) without the hypothesis that X is finitely generated. To do so, we will need axiom A4 as well as an axiom of set theory called *Zorn's lemma*. This axiom is equivalent to the Axiom of Choice, which in turn can be formulated in many equivalent ways. When applying the Axiom of Choice to other parts of mathematics, it is often most convenient to use Zorn's lemma, so that is the version of the axiom presented here. For more details on the Axiom of Choice and its variants, the reader should consult the set theory texts by Halmos [26] or Monk [39].

Before stating Zorn's lemma itself, we must recall some definitions from the theory of partially ordered sets. A *partially ordered set* (or *poset*) is a set Z and a relation \leq on Z satisfying these axioms:

PO1. *Reflexivity*: For all $x \in Z$, $x \leq x$.

PO2. *Antisymmetry*: For all $x, y \in Z$, if $x \leq y$ and $y \leq x$ then $x = y$.

PO3. *Transitivity*: For all $x, y, z \in Z$, if $x \leq y$ and $y \leq z$ then $x \leq z$.

For example, if X is any set and Z is the set of all subsets of X , then Z with the relation \subseteq (set inclusion) is a poset. More generally, we could take Z to be any collection of subsets of X and get a poset ordered by set inclusion.

For any poset (Z, \leq) , $x \in Z$ is a *maximal element* of Z iff for all $y \in Z$, $x \leq y$ implies $x = y$. This definition says that no element of the poset is strictly larger than x under the given ordering. For example, if Z is the poset of all subsets of X satisfying a fixed logical property P , then maximal elements of the poset Z are the same as maximal subsets of X satisfying P (as defined earlier).

For any poset (Z, \leq) and any subset Y of Z , $z \in Z$ is an *upper bound* for Y in the poset iff $y \leq z$ for all $y \in Y$. A subset Y of Z is a *chain* iff for all $x, y \in Y$, $x \leq y$ or $y \leq x$. *Zorn's lemma* is the following statement: If (Z, \leq) is any poset such that every chain $Y \subseteq Z$ has an upper bound in Z , then Z has a maximal element. We adopt this statement as an axiom, so it does not require proof (although, as mentioned above, it can be proved as a consequence of the Axiom of Choice).

We frequently apply Zorn's lemma to the situation where Z is a collection of subsets of some set X and \leq is set inclusion. To verify the hypothesis of Zorn's lemma in this setting, we must start with an arbitrary chain Y of elements of Z , which we can write as an indexed

set $Y = \{S_i : i \in I\}$. Note that every S_i is a subset of X . The assumption that Y is a chain means that for all $i, j \in I$, either $S_i \subseteq S_j$ or $S_j \subseteq S_i$. If Z were the collection of *all* subsets of X , we see immediately that $S = \bigcup_{i \in I} S_i$ would be an upper bound for the chain Y , since $S_i \subseteq S$ for all $i \in I$. The trouble is that Z may not consist of *all* subsets of X , so the fact that S belongs to Z requires proof. (Indeed, for some choices of Z , the union of the S_i is not in Z , so some other subset of X must be used as an upper bound for Y in Z .) In any case, once we have found an upper bound for the chain Y that does belong to Z , we can conclude via Zorn's lemma that Z has a maximal element. By definition, this is a subset M of X such that no subset of X properly containing M belongs to Z .

One other subtlety that may occur when using Zorn's lemma is the fact that the empty subset of (Z, \leq) is always a chain. Any element of Z will serve as an upper bound for this chain. Checking that the hypothesis of Zorn's lemma does hold for this special chain amounts to proving that *the set Z is nonempty*. Depending on how Z is defined, it may be a nontrivial task to verify this assertion. Once this verification has been done, one can assume (when checking the hypothesis of Zorn's lemma) that all chains being considered are nonempty.

Before using Zorn's lemma in the context of independence structures, let us give an example of a basic result in abstract algebra whose proof requires Zorn's lemma. We will prove that any nonzero commutative ring R contains a *maximal ideal*, i.e., an ideal M different from R such that there is no ideal J with $M \subsetneq J \subsetneq R$. (For the definitions of commutative ring and ideal, see §1.2 and §1.4.)

To start the proof, we need an appropriate poset (Z, \leq) . Let Z be the set of all proper ideals of R (so R itself does not belong to Z), and order Z by set inclusion. Since R is a nonzero ring, the set $\{0_R\}$ is a proper ideal of R , and hence Z is nonempty. To check the hypothesis of Zorn's lemma, let $\{J_t : t \in T\}$ be any nonempty chain of elements of Z ; so each J_t is a proper ideal of R . Our candidate for an upper bound for this chain is $J = \bigcup_{t \in T} J_t$. We must check that J does belong to Z , i.e., that J is a proper ideal of R . On one hand, 1_R does not belong to any J_t for $t \in T$ (otherwise, one readily checks that $J_t = R$, but then J_t is not a proper ideal of R). Then 1_R does not belong to the union J of the J_t 's, so that J is a proper subset of R . On the other hand, let us check that J is an ideal of R . First, $0_R \in J$ since $0_R \in J_t$ for every J_t (this uses the fact that $T \neq \emptyset$). Next, fix $x, y \in J$ and $r \in R$. We have $x \in J_s$ and $y \in J_t$ for some $s, t \in T$. Then $-x \in J_s$ (since J_s is an ideal), hence $-x \in J$. Similarly, $rx \in J_s$ and so $rx \in J$. Finally, is $x + y \in J$? Since we are looking at a chain of ideals, either $J_s \subseteq J_t$ or $J_t \subseteq J_s$. In the first case, x and y both belong to J_t , so $x + y \in J_t \subseteq J$. In the second case, $x, y \in J_s$, so $x + y \in J_s \subseteq J$. Thus, J is an ideal. We now see that J is indeed an upper bound for the given chain in the poset Z .

Having checked the hypothesis of Zorn's lemma, we conclude from the lemma that the poset Z has a maximal element. By definition, this is a proper ideal M of R such that there does not exist $N \in Z$ with $M \subsetneq N$. In other words, there is no proper ideal of R properly containing M , which is exactly what it means for M to be a maximal ideal of R .

16.7 Main Theorems: General Case

We are almost ready to give the infinite-dimensional versions of the theorems in §16.5. First we need to recall the following definitions from the theory of infinite sets. For sets S and T , we write $|S| = |T|$ iff there exists a bijection (one-to-one onto map) $f : S \rightarrow T$. We write $|S| \leq |T|$ iff there exists an injection (one-to-one map) $f : S \rightarrow T$. The *Schröder–Bernstein Theorem* asserts that $|S| = |T|$ iff $|S| \leq |T|$ and $|T| \leq |S|$. This theorem is intuitively

evident when S and T are finite sets, but it is quite tricky to prove for arbitrary infinite sets S and T ; see [55, p. 29] for a nice proof.

- T20. *For all $S \subseteq X$, S is independent iff every finite subset of S is independent.*

Proof: First assume S is independent.

Then every finite subset of S is independent (by T6).

Conversely, assume S is dependent.

Choose $x \in S$ with $x \in \text{Sp}(S - x)$ (by D3).

Choose a finite subset $F \subseteq S - x$ with $x \in \text{Sp}(F)$ (by A4).

Let $U = F + x$. U is a finite subset of S , and $U - x = F$.

We have $x \in U$ and $x \in \text{Sp}(U - x)$ (since $x \in \text{Sp}(F)$).

Thus, there exists a finite subset U of S that is dependent (by D3).

- T21. *Let $\{S_i : i \in I\}$ be a chain of independent subsets of X . Then $S = \bigcup_{i \in I} S_i$ is independent.*

Proof: Assume $\{S_i : i \in I\}$ is a chain of independent subsets of X .

To prove independence of S , we prove the independence of an arbitrary finite subset $\{y_1, \dots, y_k\}$ of S (by T20).

For $1 \leq j \leq k$, choose $i_j \in I$ with $y_j \in S_{i_j}$ (since $y_j \in S = \bigcup_{i \in I} S_i$).

For each $i, i' \in I$, either $S_i \subseteq S_{i'}$ or $S_{i'} \subseteq S_i$ (definition of a chain).

By induction on k , there exists $i_0 \in I$ with $S_{i_j} \subseteq S_{i_0}$ for $1 \leq j \leq k$.

We have $y_j \in S_{i_0}$ for $1 \leq j \leq k$ (as $y_j \in S_{i_j}$).

So $\{y_1, \dots, y_k\}$ is a subset of the independent set S_{i_0} .

Therefore $\{y_1, \dots, y_k\}$ is independent (by T6).

- T22. *Given $U \subseteq W \subseteq X$ with U independent, there exists a maximal independent subset of W that contains U .*

Proof: Let Z be the set of all independent subsets of W that contain U . We know (Z, \subseteq) is a poset. Z is nonempty, since $U \in Z$. We show that any nonempty chain in the poset Z has an upper bound in Z . Let $\{S_i : i \in I\}$ be such a chain. Then $S = \bigcup_{i \in I} S_i$ is independent (by T21). Also, $U \subseteq S \subseteq W$, since this is true of each S_i . So, S does lie in Z and is an upper bound for the given chain. Therefore, Z satisfies the hypotheses of Zorn's lemma. By that lemma, Z contains a maximal element M . The definition of Z shows that M is a maximal independent set containing U and contained in W .

- T23. *Any independent set in X can be extended to a basis of X .*

Proof: Let U be the given independent set.

Choose a maximal independent $M \subseteq X$ containing U (by T22 with $W = X$).

M is a basis of X containing U (by T13).

- T24. *There exists a basis of X .*

Proof: Apply T23, starting with the independent set \emptyset (see T1).

- T25. *Any spanning set T of X contains a basis of X .*

Proof: Choose a maximal independent subset M of T (by T22 with $U = \emptyset$ and $W = T$).

M is a basis of X contained in T (by T12).

- T26. *If S is any independent set in X and T is any spanning set for X , then $|S| \leq |T|$.*

Proof: Fix the independent set S and the spanning set T . Recalling the definition of $|S| \leq |T|$, we must construct an injection (one-to-one map) $h : S \rightarrow T$. The idea is to use Zorn's lemma to assemble “partial” injections mapping proper subsets of S into T to get an injection with the largest possible domain. If this domain is not all of S , we use T14 to make the domain even larger, which will cause a contradiction.

Step 1: We define the poset.

Let Z be the set of triples (f, A, C) where $S \cap T \subseteq A \subseteq S$, $S \cap T \subseteq C \subseteq T$, $f : A \rightarrow C$ is a bijection, $f(z) = z$ for all $z \in S \cap T$, and $(S \sim A) \cup C$ is independent. Partially order Z by defining $(f, A, C) \leq (f_1, A_1, C_1)$ iff $A \subseteq A_1$, $C \subseteq C_1$, and $f \subseteq f_1$. (Here we are viewing the functions f and f_1 as sets of ordered pairs, e.g., $f = \{(a, f(a)) : a \in A\}$. The condition $f \subseteq f_1$ means that $f_1(a) = f(a)$ for all $a \in A \subseteq A_1$, i.e., f_1 extends f .) It is routine to check the poset axioms for (Z, \leq) .

Step 2: We check the hypothesis of Zorn's lemma.

First, is the poset Z nonempty? Yes, as one checks that $(\text{id}_{S \cap T}, S \cap T, S \cap T)$ is in Z (note that $(S \sim (S \cap T)) \cup (S \cap T) = S$ is independent).

Second, given a nonempty chain $\{(f_i, A_i, C_i) : i \in I\} \subseteq Z$, we must find an upper bound (f, A, C) for this chain that lies in Z . We let $A = \bigcup_{i \in I} A_i$, $C = \bigcup_{i \in I} C_i$, and $f = \bigcup_{i \in I} f_i$. It will be evident that (f, A, C) is an upper bound of the given chain, once we show that $(f, A, C) \in Z$. For this, one must first check that f , which is a certain set of ordered pairs, is in fact a (single-valued) function with domain A . In more detail, one must check that for all $a \in A$, there exists a unique $c \in X$ with $c \in C$ and $(a, c) \in f$. Fix $a \in A$. Now $a \in A_i$ for some $i \in I$, so there is $c \in C_i \subseteq C$ with $(a, c) \in f_i \subseteq f$, namely $c = f_i(a)$. Suppose we also have $(a, d) \in f$ for some $d \in X$; we must prove $d = c$. Note $(a, d) \in f$ means $f_j(a) = d$ for some $j \in I$. We are dealing with a chain, so $(f_i, A_i, C_i) \leq (f_j, A_j, C_j)$ or $(f_j, A_j, C_j) \leq (f_i, A_i, C_i)$. In the first case, f_j extends f_i , so $d = f_j(a) = f_i(a) = c$; similarly in the other case. Analogous reasoning shows that f maps A one-to-one onto C , so that we have a well-defined bijection $f : A \rightarrow C$. Since f extends every f_i , we also have $f(z) = z$ for all $z \in S \cap T$. Note $S \cap T \subseteq A \subseteq S$ and $S \cap T \subseteq C \subseteq T$ since these set inclusions hold for every A_i and C_i . The key point still to be checked is that $(S \sim A) \cup C$ is independent.

By T20, it is enough to prove that an arbitrary finite subset F of $(S \sim A) \cup C$ is independent. Fix such an F , and write $F = \{d_1, \dots, d_m, c_1, \dots, c_n\}$ where each $d_j \in S \sim A$ and each $c_k \in C$. For each k , there is an index $i(k) \in I$ with $c_k \in C_{i(k)}$. Since we are working with a chain and n is finite, there is a single index $i_0 \in I$ with $c_k \in C_{i_0}$ for $1 \leq k \leq n$. Moreover, $d_j \in (S \sim A) \subseteq (S \sim A_{i_0})$ for $1 \leq j \leq m$. Hence, F is a finite subset of the independent set $(S \sim A_{i_0}) \cup C_{i_0}$, and so F is independent (by T6). Thus, $(S \sim A) \cup C$ is independent, completing the proof that (f, A, C) does lie in Z .

Step 3: We analyze a maximal element of Z .

By Zorn's lemma, there is a maximal element (f, A, C) in the poset Z . We claim that $A = S$. Otherwise, there exists an element $x \in S \sim A$. As A contains $S \cap T$, x is not in T , and so $x \in ((S \sim A) \cup C) \sim T$. We can now apply T14 to the independent set $(S \sim A) \cup C$ and the spanning set T to get $y \in T \sim ((S \sim A) \cup C)$ with $((S \sim A) \cup C) - x + y$ independent. Note $y \in T \sim C$. Let $A_1 = A + x$, $C_1 = C + y$, and define $f_1 : A_1 \rightarrow C_1$ by letting $f_1(z) = f(z)$ for all $z \in A$ (hence, in particular, $f_1(z) = z$ for all $z \in S \cap T$) and $f_1(x) = y$. As $x \notin A$ and $y \notin C$, it is immediate that f_1 is a bijection extending f . Finally, $(S \sim A_1) \cup C_1 = (((S \sim A) \cup C) - x) + y$ is independent. So (f_1, A_1, C_1) is an element of the poset Z strictly larger than (f, A, C) , which contradicts maximality of (f, A, C) . So $A = S$ after all. Finally, by enlarging the codomain, we can regard the bijection $f : A \rightarrow C$ (where $C \subseteq T$) as an injection of $S = A$ into T . Therefore, $|S| \leq |T|$.

T27. If A and B are two bases of X , then $|A| = |B|$.

Proof: Since A is independent and B spans X , $|A| \leq |B|$ by T26.

Since B is independent and A spans X , $|B| \leq |A|$ by T26.

Therefore, $|A| = |B|$ (equality of cardinals) by the Schröder–Bernstein Theorem.

Having proved T24 and T27, it now makes sense to define the *dimension* of X to be the cardinality of any basis of X .

16.8 Bases of Subspaces

Let X be a set with spanning operator $\text{Sp} = \text{Sp}_X$. Let V be a subspace of X , so that V is a subset of X of the form $V = \text{Sp}(S)$ for some $S \subseteq X$. Consider the operator Sp_V defined on the set of all subsets of V by the rule $\text{Sp}_V(T) = \text{Sp}_X(T)$ for all $T \subseteq V$. For any $T \subseteq V$, we see that $\text{Sp}_V(T) = \text{Sp}_X(T) \subseteq \text{Sp}_X(V) = \text{Sp}_X(\text{Sp}_X(S)) = \text{Sp}_X(S) = V$ by T3 and T4. This proves that the pair (V, Sp_V) satisfies axiom A0. Axioms A1 through A4 automatically hold for (V, Sp_V) , since the conditions imposed by these axioms are a subset of the conditions that already hold for (X, Sp_X) , which is known to satisfy A1 through A4.

It follows that all the theorems proved for X are applicable also to all subspaces V of X . In particular, every subspace has a basis; an independent subset of a subspace can be enlarged to a basis of that subspace; a spanning set for a subspace can be shrunk to a basis of that subspace; and any two bases of a given subspace have the same cardinality. We define the *dimension* of a subspace to be the dimension of any basis of that subspace.

16.9 Linear Independence and Linear Bases

We now apply our axiomatic framework to linear algebra. Let X be a fixed vector space over a field F . (Fields are defined in §1.2. We remark that none of the definitions and results in this section will invoke commutativity of multiplication in the field F . So the theorems obtained here apply more generally to *left vector spaces* over *division rings*, as discussed in Chapter 17 below.) We begin by recalling some definitions from “classical” linear algebra.

- L0. For $S \subseteq X$, let $\text{span}_F(S)$ be the set of all finite F -linear combinations of elements of S , i.e.,

$$\text{span}_F(S) = \{0\} \cup \{c_1 s_1 + \cdots + c_n s_n : c_i \in F, s_i \in S, n \in \mathbb{N}^+\}.$$

- L1. For $S, V \subseteq X$, we say that S spans V (or generates V) iff $V = \text{span}_F(S)$.
- L2. For $V \subseteq X$, V is a subspace iff $0 \in V$, and $v + w \in V$ for all $v, w \in V$, and $cv \in V$ for all $c \in F$ and $v \in V$.
- L3. For $S \subseteq X$, S is linearly dependent iff there exist an integer $n > 0$, scalars $c_1, \dots, c_n \in F$ that are not all zero, and distinct vectors $x_1, \dots, x_n \in S$ such that $c_1 x_1 + \cdots + c_n x_n = 0$.
- L4. For $S \subseteq X$, S is linearly independent iff S is not linearly dependent.
- L5. For $S \subseteq X$, S is a (linear) basis for X iff S is independent and $\text{span}_F(S) = X$.
- L6. X is finitely generated (or finite-dimensional) iff $X = \text{span}_F(S)$ for some finite set $S \subseteq X$.

Define $\text{Sp}(S) = \text{span}_F(S)$ for all $S \subseteq X$. To apply our general theory, we must verify axioms A0 through A4.

Proof of A0: For $S \subseteq X$, we have $\text{Sp}(S) \subseteq X$, since X is closed under the vector space operations.

Proof of A1: Suppose $S \subseteq X$ and $v \in S$. Then $v = 1v$ is a linear combination of elements of S , so $v \in \text{Sp}(S)$ and $S \subseteq \text{Sp}(S)$.

Proof of A2: Suppose $S, T \subseteq X$ and $S \subseteq \text{Sp}(T)$. We show that $\text{Sp}(S) \subseteq \text{Sp}(T)$. Let $v \in \text{Sp}(S)$. Since $0 \in \text{Sp}(T)$, we may assume $v \neq 0$ and write $v = c_1x_1 + \cdots + c_nx_n$ for some $n > 0$, $c_i \in F$, and $x_i \in S$. Because $S \subseteq \text{Sp}(T)$, we can write each x_i as a linear combination of elements of T , say

$$x_i = \sum_{j=1}^m b_{i,j}y_j$$

for some $m < \infty$, $b_{i,j} \in F$, and $y_j \in T$. (By allowing zero coefficients, we can assume the same set of y_j 's appears in the expansion of every x_i .) Then

$$v = \sum_{i=1}^n c_i x_i = \sum_{i=1}^n c_i \sum_{j=1}^m b_{i,j} y_j = \sum_{j=1}^m \left(\sum_{i=1}^n c_i b_{i,j} \right) y_j,$$

where each coefficient $\sum_{i=1}^n c_i b_{i,j}$ is in F . This shows that v is a linear combination of elements of T . So $v \in \text{Sp}(T)$.

Proof of A3: Suppose $x, y \in X$, $S \subseteq X$, $x \in \text{Sp}(S + y)$, and $x \notin \text{Sp}(S)$. We show $y \in \text{Sp}(S + x)$. First, write

$$x = by + c_1x_1 + \cdots + c_nx_n$$

where b and each c_i are scalars in F , and each $x_i \in S$. Note that $b \neq 0$, since otherwise the displayed relation implies that $x \in \text{Sp}(S)$. Hence, b^{-1} exists in the field F , and solving for y leads to

$$y = b^{-1}x - b^{-1}c_1x_1 - \cdots - b^{-1}c_nx_n.$$

This shows that $y \in \text{Sp}(S + x)$.

Proof of A4: Suppose $S \subseteq X$, $x \in X$, and $x \in \text{Sp}(S)$. If $x = 0$, then $x \in \text{Sp}(\emptyset)$, where \emptyset is a finite subset of S . Otherwise, by definition of $\text{span}_F(S)$, x is a *finite* linear combination of elements of S , say

$$x = c_1x_1 + \cdots + c_nx_n$$

where $c_i \in F$, $x_i \in S$, and $n < \infty$. Putting $S' = \{x_1, \dots, x_n\}$, S' is a finite subset of S such that $x \in \text{Sp}(S')$.

Next, we check that the classical definitions L1 through L6 are equivalent to the corresponding formal definitions D1 through D6.

(D1 \Leftrightarrow L1): This equivalence is immediate, since $\text{Sp}(S) = \text{span}_F(S)$.

(D2 \Leftrightarrow L2): Assume V is a subspace of X as in D2, say $V = \text{Sp}(S)$ for some $S \subseteq X$. We have $0 \in V$, by definition of $\text{span}_F(S)$. Let $v, w \in V$ and $b \in F$. Assuming (as we may) that $v \neq 0 \neq w$, we can write $v = c_1x_1 + \cdots + c_nx_n$ and $w = d_1x_1 + \cdots + d_nx_n$, where $c_i, d_i \in F$ and $x_i \in S$. Then

$$v + w = (c_1 + d_1)x_1 + \cdots + (c_n + d_n)x_n \in V,$$

$$bv = (bc_1)x_1 + \cdots + (bc_n)x_n \in V,$$

and so V is a linear subspace of X as in L2.

Conversely, assume V is a linear subspace as in L2. By A1, $V \subseteq \text{Sp}(V)$. For the reverse inclusion, take $c_1v_1 + \cdots + c_nv_n \in \text{Sp}(V)$, where $c_i \in F$ and $v_i \in V$. Since V is a linear

subspace, each $c_i v_i$ lies in V . By induction on n , the sum of these vectors also lies in V . Therefore, $\text{Sp}(V) \subseteq V$, so that $V = \text{Sp}(V)$. Letting $S = V$, we have $V = \text{Sp}(S)$, so that V is a subspace as in D2.

(D3 \Leftrightarrow L3): Assume $S \subseteq X$ is dependent as defined in D3. Choose $x \in S$ satisfying $x \in \text{Sp}(S - x)$. We can write

$$x = c_1 x_1 + \cdots + c_n x_n,$$

where $c_i \in F$ and the x_i 's are distinct elements of $S - x$. Rearranging gives

$$c_1 x_1 + \cdots + c_n x_n - 1x = 0,$$

where the last coefficient is nonzero. This relation shows that S is linearly dependent as defined in L3. Conversely, if S is linearly dependent as defined in L3, there is a relation

$$d_1 y_1 + d_2 y_2 + \cdots + d_m y_m = 0,$$

where y_i are distinct elements of S , $d_i \in F$, and some $d_i \neq 0$. We may choose notation so that $d_1 \neq 0$. Then

$$y_1 = -d_1^{-1} d_2 y_2 - \cdots - d_1^{-1} d_m y_m.$$

Thus $y_1 \in \text{Sp}(S - y_1)$, so S is dependent as defined in D3.

(D4 \Leftrightarrow L4), (D5 \Leftrightarrow L5), and (D6 \Leftrightarrow L6): These follow immediately from the previous results.

Now we may deduce the following classical linear algebra theorems: every vector space X has a linear basis; every linearly independent subset of X can be enlarged to a linear basis; every spanning set for X contains a linear basis; linearly independent sets are never larger than spanning sets; and any two linear bases of X have the same cardinality. We define the (*linear*) dimension of X to be the cardinality of any linear basis of X . We stress that our proofs are valid even for infinite-dimensional vector spaces.

16.10 Field Extensions

In Chapter 12, we developed some facts about field extensions that helped us prove results about ruler and compass constructions. That chapter focused mainly on *finite-degree* field extensions, which consist of a field K and a subfield F such that $[K : F]$, the dimension of K viewed as an F -vector space, is finite.

For more advanced work in field theory, one would like to have a more detailed understanding of field extensions $F \subseteq K$ in which $[K : F] = \infty$. The first step in this direction is to introduce the *transcendence degree* of a field extension, which (informally) measures the extent of infinitude of an infinite-degree field extension. We can deduce the basic facts about this new concept from the axiomatic framework in this chapter. First, however, we need a few more definitions and results from field theory. Some of these results were proved in §12.3 and §12.7; proofs of the others are sketched in the exercises of this chapter.

Field extension generated by a set. Suppose A is any subset of a field X . A *monomial in A* is a finite product $a_1 \cdots a_n$, where the a_i are (not necessarily distinct) elements of A . An empty product is interpreted as 1_X . A *polynomial expression in A* is a finite sum of monomials in A ; an empty sum is interpreted as 0_X . A *rational expression in A* is an element $f/g \in X$, where f and g are polynomial expressions in A such that $g \neq 0_X$. If K is a subfield of

X and S is any subset of X , define $K(S)$ to be the set of all elements of X that can be written as rational expressions in $K \cup S$. Aided by the algebraic rules for computing with polynomials and fractions, one shows that: $K(S)$ is a subfield of X that contains K and S ; and $K(S)$ is contained in every subfield of X that contains K and S . $K(S)$ is called the *subfield generated by S over K* . Some authors define $K(S)$ as the intersection of all subfields of X containing $K \cup S$; our definition has the advantage of giving an explicit description of the form of each element of $K(S)$. When S is a one-element set $\{z\}$, we write $K(z)$ instead of $K(S)$. One can check that for $S, T \subseteq X$, $K(S \cup T) = (K(S))(T) = (K(T))(S)$.

Finite extensions. If K_1 is a subfield of K_2 , we can view K_2 as a vector space over the field K_1 . As mentioned above, $[K_2 : K_1]$ denotes the dimension of this vector space, and K_2 is a *finite extension* of K_1 iff $[K_2 : K_1]$ is finite. If K_1 is a subfield of K_2 and K_2 is subfield of K_3 , then we have the *product formula* $[K_3 : K_1] = [K_3 : K_2] \cdot [K_2 : K_1]$. The formula is valid even for dimensions that are infinite cardinals; the finite case was proved in §12.3. In particular, given subfields $K_1 \subseteq K_2 \subseteq \cdots \subseteq K_m \subseteq X$ such that K_{i+1} is a finite extension of K_i for $1 \leq i < m \in \mathbb{N}^+$, iteration of the product formula shows that K_m is a finite extension of K_1 .

Algebraic elements and extensions. Let K be a subfield of a field X , and suppose $z \in X$. We say z is *algebraic over K* iff there exists a nonzero polynomial f with coefficients in K that has z as a root. In other words, for some $n > 0$ and some $b_0, b_1, \dots, b_{n-1} \in K$, the identity $z^n + b_{n-1}z^{n-1} + \cdots + b_1z + b_0 = 0$ holds in X . We say z is *transcendental over K* iff z is not algebraic over K . A subfield E of X containing K is *algebraic over K* or an *algebraic extension of K* iff every $z \in E$ is algebraic over K .

We list without proof some facts about algebraic elements and extensions. Let X be a field with subfields $K \subseteq E \subseteq L$. First, every $z \in K$ is algebraic over K . Second, if $z \in X$ is algebraic over K , then z is algebraic over E . Third, an element $x \in X$ is algebraic over K iff the dimension $[K(x) : K]$ is finite (see §12.7). Fourth, if E is a finite extension of K , then E is an algebraic extension of K . Fifth, if E is an algebraic extension of K and L is an algebraic extension of E , then L is an algebraic extension of K . Sixth, the set of all elements $z \in X$ that are algebraic over K is a subfield of X containing K .

16.11 Algebraic Independence and Transcendence Bases

To begin our discussion of transcendence degree, let X be a field with subfield F ; F and X will remain fixed throughout the following discussion. We introduce the following field-theoretic definitions.

- F1. For $S, V \subseteq X$, S *transcendentally spans V over F* (or *transcendentally generates V over F*) iff V consists of all $x \in X$ that are algebraic over $F(S)$.
- F2. For $V \subseteq X$, V is *algebraically closed in X* iff V is a subfield of X such that $F \subseteq V$ and whenever $z \in X$ is algebraic over V , we have $z \in V$.
- F3. For $S \subseteq X$, S is *algebraically dependent over F* iff there exists $n \in \mathbb{N}^+$ and a nonzero polynomial $f \in F[x_1, \dots, x_n]$ and distinct elements $s_1, \dots, s_n \in S$ with $f(s_1, \dots, s_n) = 0$. Informally, this means that there is a nontrivial *algebraic* combination of elements of S , not merely an *F -linear* combination, that evaluates to zero. Here, “algebraic” means that we can multiply elements of S by one another, as well as multiplying them by scalars from F and adding together the resulting quantities.

- F4. For $S \subseteq X$, S is *algebraically independent over F* iff S is not algebraically dependent over F . In more detail, this means that for all $n \in \mathbb{N}^+$, all nonzero $f \in F[x_1, \dots, x_n]$, and all distinct $s_1, \dots, s_n \in S$, $f(s_1, \dots, s_n) \neq 0_X$.
- F5. For $S \subseteq X$, S is a *transcendence basis* for X over F iff S is algebraically independent over F and S transcendentally spans X .
- F6. X has *finite transcendence degree* over F iff X is transcendentally spanned by some finite set $S \subseteq X$.

We now apply our axiomatic setup to the current situation. For $S \subseteq X$, define $\text{Sp}(S)$ to be the set of all $x \in X$ that are algebraic over $F(S)$. We must prove axioms A0 through A4.

Proof of A0: For $S \subseteq X$, the definition just given shows that $\text{Sp}(S) \subseteq X$.

Proof of A1: We show that $S \subseteq \text{Sp}(S)$ for all $S \subseteq X$. Given $s \in S$, s belongs to $F(S)$, so s is algebraic over $F(S)$. Hence, $s \in \text{Sp}(S)$.

Proof of A4: Suppose $z \in \text{Sp}(S)$, so that z is algebraic over $F(S)$. This means that there is a polynomial $f = x^n + b_{n-1}x^{n-1} + \dots + b_1x + b_0$ with each $b_i \in F(S)$ and $f(z) = 0$. Every b_i is a rational expression in $F \cup S$; by definition, such an expression is built up by multiplying, adding, and dividing a finite number of elements from $F \cup S$. Also, there are only finitely many b_i 's. Therefore, we can find a finite subset S' of S (consisting of all elements of S appearing in the rational expressions for the b_i 's) such that every b_i is a rational expression in $F \cup S'$. Then each b_i lies in $F(S')$, so that the coefficients of f lie in $F(S')$. Since $f(z) = 0$, this shows that z is algebraic over $F(S')$. Therefore, $z \in \text{Sp}(S')$ for some finite subset S' of S .

Proof of A2: Let $S, T \subseteq X$ satisfy $S \subseteq \text{Sp}(T)$. We must show that $\text{Sp}(S) \subseteq \text{Sp}(T)$. Let $z \in \text{Sp}(S)$. By A4, choose a finite subset $S' = \{s_1, \dots, s_m\}$ of S with $z \in \text{Sp}(S')$. Define fields $F_0 = F(T)$, $F_i = F(T \cup \{s_1, \dots, s_i\})$ for $1 \leq i \leq m$, and $F_{m+1} = F(T \cup S' \cup \{z\})$. Then we have a chain of field extensions

$$F_0 \subseteq F_1 \subseteq \dots \subseteq F_m \subseteq F_{m+1}.$$

By the assumption $S \subseteq \text{Sp}(T)$, each s_i is algebraic over $F(T) = F_0$, and therefore s_i is algebraic over the larger field F_{i-1} . Similarly, z is algebraic over $F(S')$, so z is algebraic over the larger field $F_m = F(T \cup S')$. Since $F_i = F_{i-1}(s_i)$ for $1 \leq i \leq m$ and $F_{m+1} = F_m(z)$, it follows that each dimension $[F_i : F_{i-1}]$ is finite. By the product formula, $[F_{m+1} : F_0] = \prod_{i=1}^{m+1} [F_i : F_{i-1}]$ is also finite. Since $z \in F_{m+1}$, z must be algebraic over $F_0 = F(T)$, and so $z \in \text{Sp}(T)$.

Proof of A3: Assume $u, v \in X$ and $S \subseteq X$ satisfy $v \in \text{Sp}(S+u)$ but $v \notin \text{Sp}(S)$; we must prove that $u \in \text{Sp}(S+v)$. The idea of the proof is to construct a certain two-variable polynomial $h = h(x, y) \in K[x, y]$, where K is the field $F(S)$ and x and y are formal variables. By evaluating x at u , h becomes a one-variable polynomial $h(u, y) \in K(u)[y] = F(S+u)[y]$. On the other hand, by evaluating y at v , h becomes a one-variable polynomial $h(x, v) \in K(v)[x] = F(S+v)[x]$. The key to the proof will be transferring information between these two specializations of h . (For more information on formal polynomials and their evaluations, see Chapter 3, especially §3.5 and §3.20.)

To begin, the assumption $v \in \text{Sp}(S+u)$ means that v is algebraic over $F(S+u)$. This means there is a nonzero polynomial $f = \sum_{i=0}^n b_i y^i \in F(S+u)[y]$ with $f(v) = 0$ and each $b_i \in F(S+u)$. Each b_i is a rational expression in $(F \cup S) \cup \{u\}$, say $b_i = g_i/k_i$ for certain polynomial expressions g_i, k_i in $(F \cup S) \cup \{u\}$. Multiplying both sides of $f(v) = 0$ by $k_0 k_1 \dots k_n$, we can eliminate all the denominators in the b_i 's. Changing notation if needed, we can now assume $f(v) = 0$ where each b_i is a polynomial expression in $(F \cup S) \cup \{u\}$.

Then each b_i has the form $b_i = \sum_{j \geq 0} a_{i,j} u^j$ for some $a_{i,j} \in F(S)$ (in fact, each $a_{i,j}$ is a polynomial expression in $F \cup S$). Now, introduce the two-variable polynomial

$$h = h(x, y) = \sum_{i=0}^n \sum_{j \geq 0} a_{i,j} x^j y^i \in F(S)[x, y].$$

By construction, $h(u, y) = \sum_{i=0}^n (\sum_{j \geq 0} a_{i,j} u^j) y^i = \sum_{i=0}^n b_i y^i = f$ is a nonzero polynomial in y having v as a root. So $h(u, v) = 0$.

Note that h is a nonzero polynomial, since its specialization $h(u, y)$ is nonzero. So not all $a_{i,j}$'s are zero. Let k be the maximum value of j such that $a_{i,k} \neq 0$ for some i . Can k be zero? If so, we would have $h = \sum_{i=0}^n a_{i,0} y^i = h(u, y)$ where each $a_{i,0} \in F(S)$. Since v is a root of this polynomial, we would have v algebraic over $F(S)$, contradicting the assumption $v \notin \text{Sp}(S)$. So $k > 0$. Now, evaluate $h(x, y)$ at $y = v$ to get

$$h(x, v) = \sum_{j=0}^k \left(\sum_{i \geq 0} a_{i,j} v^i \right) x^j.$$

This is a polynomial in $F(S + v)[x]$, which we claim is nonzero. To see why, note that the leading coefficient of $h(x, v)$ (i.e., the coefficient of x^k) is $\sum_{i \geq 0} a_{i,k} v^i$. If this were zero, then v would be a root of the nonzero polynomial $\sum_{i \geq 0} a_{i,k} y^i \in F(S)[y]$, implying that v is algebraic over $F(S)$. As above, this contradicts $v \notin \text{Sp}(S)$. Summarizing, we have found a nonzero polynomial $h(x, v) \in F(S + v)[x]$ that has u as a root, since $h(u, v) = 0$. This means u is algebraic over $F(S + v)$, so that $u \in \text{Sp}(S + v)$, as we needed to show.

Next we check that the field-theoretic definitions F1 through F6 are equivalent to the corresponding definitions D1 through D6 in the axiomatic framework.

(F1 \Leftrightarrow D1): This is immediate from the way we defined the spanning operator.

(F2 \Leftrightarrow D2): Assume V is algebraically closed in X , as defined in F2. We claim that $V = \text{Sp}(S)$ for $S = V$, so that V is a subspace as defined in D2. We have $V \subseteq \text{Sp}(V)$ by A1. On the other hand, $F(V) = V$ since the intersection of all subfields of X containing $F \cup V$ is V (as V itself is one of these subfields). If $z \in X$ is algebraic over $F(V) = V$, then $z \in V$ since V is algebraically closed in X . Hence, $\text{Sp}(V) \subseteq V$.

Conversely, assume V is a subspace as defined in D2, say $V = \text{Sp}(S)$ for some $S \subseteq X$. V is the set of elements in X that are algebraic over $F(S)$, which is a subfield of X containing F . Furthermore, if $z \in X$ is algebraic over V , then we have a chain of field extensions $F(S) \subseteq V \subseteq V(z)$. We know V is an algebraic extension of $F(S)$; moreover, $V(z)$ is a finite extension of V , hence is an algebraic extension of V . Then $V(z)$ is an algebraic extension of $F(S)$. So z is algebraic over $F(S)$, which gives $z \in \text{Sp}(S) = V$. This proves that V is algebraically closed in X .

(F3 \Leftrightarrow D3): Assume $S \subseteq X$ is dependent as defined in D3. Choose $u \in S$ with $u \in \text{Sp}(S - u)$. By A4, choose a finite $S' = \{v_1, \dots, v_m\} \subseteq S - u$ with $u \in \text{Sp}(S')$. Then u is algebraic over $F(S')$. As in the proof of A3, we can find a nonzero polynomial $f = \sum_{i=0}^n b_i x^i \in F(S')[x]$ with $f(u) = 0$ and with each b_i a polynomial expression in $F \cup S'$. Writing out what this means, we obtain a nonzero multivariable polynomial $h \in F[x, y_1, \dots, y_m]$ such that $h(x, v_1, \dots, v_m) = f$ and $h(u, v_1, \dots, v_m) = f(u) = 0$. Since u, v_1, \dots, v_m are distinct elements of S , this proves the algebraic dependence of S as defined in F3.

Conversely, assume S is algebraically dependent over F , as in F3. There is a nonzero polynomial $f \in F[x_1, \dots, x_n]$ and distinct $s_1, \dots, s_n \in S$ with $f(s_1, \dots, s_n) = 0$. Among all such dependence relations, assume we have chosen one with n minimal. Now f is not

a constant polynomial, so some variable (say x_1) occurs in a monomial of f that has a nonzero coefficient. Consider $g(x_1) = f(x_1, s_2, \dots, s_n)$, which is a one-variable polynomial with coefficients in $F(S - s_1)$ having s_1 as a root. By minimality of n and choice of x_1 , the leading coefficient of g cannot be zero. So $g \neq 0$, and s_1 is therefore algebraic over $F(S - s_1)$. This means $s_1 \in \text{Sp}(S - s_1)$, so that S is dependent as defined in D3.

(F4 \Leftrightarrow D4), (F5 \Leftrightarrow D5), and (F6 \Leftrightarrow D6): These follow immediately from the previous results.

With no further work, we can now harvest the main theorems on transcendental field extensions: every field extension X of F has a transcendence basis over F ; every algebraically independent subset of X over F can be enlarged to a transcendence basis of X ; every transcendent spanning set of X over F contains a transcendence basis of X ; algebraically independent sets are never larger than transcendent spanning sets; and any two transcendence bases of X have the same cardinality. We define the *transcendence degree* of X over F to be the cardinality of any transcendence basis of X over F .

16.12 Independence in Graphs

Our third application involves spanning forests for graphs. A *graph* G is a triple (V, E, ϵ) , where V is a set of *vertices* (possibly infinite), E is a collection of *edges*, and ϵ is an *endpoint function* with domain E such that, for each $e \in E$, $\epsilon(e)$ is a subset of V of size 1 or 2. If $\epsilon(e) = \{v\}$, then e is a *loop edge* at the vertex v ; if $\epsilon(e) = \{v, w\}$ with $v \neq w$, e is an edge with *endpoints* v and w . Whenever $v \in \epsilon(e)$, we say the vertex v and the edge e are *incident* to each other. Finally, edges $e, f \in E$ with $\epsilon(e) = \epsilon(f)$ are called *parallel edges*. A *simple graph* is a graph with no loop edges and no parallel edges e, f with $e \neq f$.

Given a graph $G = (V, E, \epsilon)$ and $v, w \in V$, a *walk* in G from v to w is a *finite* sequence $v = v_0, e_1, v_1, e_2, v_2, \dots, e_n, v_n = w$ such that $n \in \mathbb{N}$, each $v_i \in V$, each $e_i \in E$, and $\epsilon(e_i) = \{v_{i-1}, v_i\}$ for $1 \leq i \leq n$. A *trail* is a walk in which no edge is used twice; a *path* is a trail in which v_0, \dots, v_n are distinct; a *cycle* is a trail in which $n > 0$, $v_0 = v_n$, and v_0, \dots, v_{n-1} are distinct. Given $v, w \in V$, there is a walk from v to w using edges in some subset $E' \subseteq E$ iff there is a trail from v to w using edges in E' iff there is a path from v to w using edges in E' iff there is a path from w to v using edges in E' . We leave the rigorous proof of this assertion as an exercise; the idea of the proof is to keep removing redundant edge traversals to go from a walk to a trail, and then remove redundant cycles to go from a trail to a path.

Now fix a graph $G = (V, E, \epsilon)$ with no loop edges. Let $X = E$, the edge set of G . We define an independence structure on X as follows. For $S \subseteq X$, define $\text{Sp}(S)$ to be the set of edges $e \in E$ with endpoints v, w such that there exists a path from v to w using only edges in S . Axiom A0 certainly holds. If $e \in S$ has endpoints v, w , then the sequence v, e, w is a path from v to w using the edge $e \in S$, so that $e \in \text{Sp}(S)$. Thus, $S \subseteq \text{Sp}(S)$, and axiom A1 holds. For axiom A2, we assume that $S \subseteq \text{Sp}(T)$ and prove that $\text{Sp}(S) \subseteq \text{Sp}(T)$. Given $e \in \text{Sp}(S)$ with endpoints v, w , let $v = v_0, e_1, v_1, \dots, e_n, v_n = w$ be a path such that $e_i \in S$ for $1 \leq i \leq n$. Since $S \subseteq \text{Sp}(T)$, we can replace each subpath v_{i-1}, e_i, v_i by a path from v_{i-1} to v_i with all its edges in T . Concatenating these new paths, we obtain a walk from $v = v_0$ to $v_n = w$ consisting of edges in T . Shrinking this walk to a path, we see that $e \in \text{Sp}(T)$, as needed.

For axiom A3, assume $e \in \text{Sp}(S + f)$ and $e \notin \text{Sp}(S)$, where $\epsilon(e) = \{v, w\}$ and $\epsilon(f) = \{x, y\}$. We must show that $f \in \text{Sp}(S + e)$. There is a path $v_0, e_1, v_1, \dots, e_n, v_n$ from v to

w with each $e_j \in S + f$. One of these edges, say e_i , must be the edge f , since otherwise e would be in $\text{Sp}(S)$. Since we are dealing with a path, this is the only occurrence of edge f in the given path. Choose notation so that $x = v_{i-1}$ and $y = v_i$. Then

$$x = v_{i-1}, e_{i-1}, v_{i-2}, \dots, e_1, v_0 = v, e, w = v_n, e_n, v_{n-1}, \dots, e_{i+1}, v_i = y$$

is a path from x to y using edges in $S + e$. So $f \in \text{Sp}(S + e)$, and axiom A3 holds.

Finally, to verify axiom A4, assume $S \subseteq X$ and $e \in X$ satisfy $e \in \text{Sp}(S)$. Letting $\epsilon(e) = \{v, w\}$, we know there is a path $v_0, e_1, v_1, \dots, e_n, v_n$ from v to w . By definition, this path can use only a finite number of edges from S . Letting S' be the finite set of edges used, we then have $e \in \text{Sp}(S')$.

Now, $S \subseteq E$ is dependent iff $e \in \text{Sp}(S - e)$ for some $e \in S$ iff there exists $e \in S$ (with endpoints v, w) and a path from v to w using only edges in S different from e . Appending w, e, v to this path, we see that dependence of S is equivalent to the existence of a cycle in G using edges in S . S is independent iff there are no such cycles, which we abbreviate by saying S is an *acyclic* edge set. S is a basis iff S is a maximal independent set (by T7 and T13). One can check that maximality of S implies that every vertex of G that is incident to at least one edge of E must also be incident to at least one edge of S . Accordingly, bases in this setting are also called *spanning forests* for G .

We conclude from the general theory that every graph has a spanning forest; the cardinality of a spanning forest of G is unique (even for G infinite); and every acyclic edge set in G can be enlarged to a spanning forest by adding appropriate edges.

16.13 Hereditary Systems

We conclude this chapter with a very brief introduction to the vast subject of matroids. We mentioned in §16.1 that a *matroid* is a pair (X, Sp) , where X is a *finite* set and the spanning operator satisfies axioms A0, A1, A2, and A3 (axiom A4 follows from the finiteness of X). However, there are many other equivalent ways of defining matroids. Before describing these alternate definitions, we study the more general idea of a hereditary system.

Given any poset (Z, \leq) , an *order-ideal* of Z is a nonempty subset \mathcal{I} of Z such that for all $u, v \in Z$, if $u \leq v$ and $v \in \mathcal{I}$, then $u \in \mathcal{I}$. We will focus exclusively on the special case where X is a finite set and Z is the poset of all subsets of X ordered by set inclusion. In this case, an order-ideal of Z is called a *hereditary system with ground set X* . Restating the definition in this special case, we see that a hereditary system on X is a nonempty collection \mathcal{I} of subsets of X such that every subset of a set in \mathcal{I} is also in \mathcal{I} . Given that \mathcal{I} is “closed under taking subsets,” one sees that the condition $\mathcal{I} \neq \emptyset$ is equivalent to the condition $\emptyset \in \mathcal{I}$. The subsets in \mathcal{I} are called the *independent sets* of the hereditary system. For example, theorems T1 and T6 show that the independent sets in an independence structure (X, Sp) (as defined in D4) form a hereditary system with ground set X .

Now let \mathcal{I} be an arbitrary hereditary system on a finite set X . By analogy with the examples discussed earlier, we can use \mathcal{I} to define various auxiliary concepts:

- H1. *Dependent Sets.* A subset S of X is called *dependent* iff $S \notin \mathcal{I}$ iff S is not independent. Note that \emptyset is not dependent, and any superset of a dependent set is also dependent.
- H2. *Bases.* Define a *basis* of the hereditary system to be a maximal element of \mathcal{I} . So, $S \subseteq X$ is a basis iff S is independent and no set $T \subseteq X$ properly containing S is

independent. Equivalently, one checks that S is a basis iff S is independent and $S + x$ is dependent for all $x \in X \sim S$.

- H3. *Circuits.* A subset C of X is called a *circuit* iff C is a minimal dependent subset of X . So, C is a circuit iff C is dependent and no proper subset of C is dependent. Equivalently, one sees that C is a circuit iff C is dependent and $C - x$ is independent for all $x \in C$. (The terminology “circuit” comes from applications to graph theory; cf. §16.12.)
- H4. *Spanning Sets.* A subset S of X is called a *spanning set* iff S contains some basis of X .
- H5. *Rank Function.* Define a *rank function* on the set of all subsets $S \subseteq X$ by letting $\text{rk}(S)$ be the maximum size of any independent subset of S .
- H6. *Spanning Operator.* Define a *spanning operator* on the set of all subsets $S \subseteq X$ by letting

$$\text{Sp}(S) = S \cup \{z \in X : \text{for some } T \subseteq S, T + z \text{ is a circuit.}\} \quad (16.1)$$

We see from these definitions that knowing the independent sets (members of \mathcal{I}) completely determines the dependent sets, bases, circuits, spanning sets, rank function, and spanning operator for X . For example, we can find all the circuits (in principle) by listing all subsets of X not in \mathcal{I} and picking out the minimal elements of this collection of sets relative to set inclusion. This, in turn, allows us to use (16.1) to compute $\text{Sp}(S)$ for any subset S of X . Conversely, one can show that any of the six concepts defined above could have been used as the starting point for defining a hereditary system. Each concept obeys a short list of axioms, and anything satisfying those axioms determines a unique hereditary system \mathcal{I} . For instance, the collection \mathcal{B} of bases of a hereditary system \mathcal{I} on X is nonempty and has the property that for all $B, C \in \mathcal{B}$, B is not a proper subset of C . Given any collection \mathcal{B}' of subsets of X satisfying this condition, one can show there exists a unique hereditary system \mathcal{I} on X having \mathcal{B}' as its set of bases. There are similar results that allow us to define a hereditary system by specifying its dependent sets or circuits. One can also give axioms characterizing the rank functions or spanning operators of hereditary systems, but these are somewhat more subtle and will not be discussed here.

16.14 Matroids

Matroids can be defined in many equivalent ways. Most commonly, one starts with a finite set X and a hereditary system \mathcal{I} with ground set X . Not all hereditary systems are matroids. To obtain a matroid, we must impose one additional “regularity condition” that forces many other nice properties to hold automatically. Several of these properties of hereditary systems are logically equivalent to one another, so we could take any one of them as the extra condition in the definition of a matroid. In a given application, we obtain all of the equivalent properties for free once we check that one of the properties is true.

We now list thirteen properties, each of which can be used as the final axiom in the definition of a matroid. The next section and the exercises contain proofs of the logical equivalence of some of these properties; for proofs of the remaining equivalences, we refer the reader to [62, Sec. 8.2] or other texts on matroids. We assume throughout that \mathcal{I} is a hereditary system on the finite set X , and that dependent sets, bases, circuits, spanning sets, the rank function, and the spanning operator are defined from \mathcal{I} as in §16.13.

- M1. *Augmentation of Independent Sets.* For all independent subsets S, T in X with $|T| > |S|$, there exists $z \in T \sim S$ such that $S + z$ is independent.
- M2. *Uniform Size of Bases.* For all $S \subseteq X$, any two maximal independent subsets of S have the same size.
- M3. *Subtract-and-Add Basis Exchange.* For all bases S, T of X and all $s \in S \sim T$, there exists $t \in T \sim S$ such that $(S - s) + t$ is a basis of X .
- M4. *Add-and-Subtract Basis Exchange.* For all bases S, T of X and all $t \in T \sim S$, there exists $s \in S \sim T$ such that $(S + t) - s$ is a basis of X .
- M5. *Submodularity of Rank Function.* For all $S, T \subseteq X$,

$$\text{rk}(S \cap T) + \text{rk}(S \cup T) \leq \text{rk}(S) + \text{rk}(T).$$

- M6. *Weak Absorption of Rank Function.* For all $S \subseteq X$ and all $u, v \in X$, if $\text{rk}(S) = \text{rk}(S + u) = \text{rk}(S + v)$ then $\text{rk}((S + u) + v) = \text{rk}(S)$.
- M7. *Strong Absorption of Rank Function.* For all $S, T \subseteq X$, if $\text{rk}(S + t) = \text{rk}(S)$ for all $t \in T$, then $\text{rk}(S) = \text{rk}(S \cup T)$.
- M8. *Weak Elimination for Circuits.* For all circuits $C \neq D$ in X and all $z \in C \cap D$, there is a circuit $C' \subseteq (C \cup D) - z$ (equivalently, $(C \cup D) - z$ is dependent).
- M9. *Strong Elimination for Circuits.* For all circuits C, D in X and all $z \in C \cap D$ and all $v \in C \sim D$, there is a circuit C' with $v \in C' \subseteq (C \cup D) - z$.
- M10. *Uniqueness of New Circuits.* For all independent sets S in X and all $z \in X$, $S + z$ contains at most one circuit.
- M11. *Rank-Preservation of Spanning Operator.* For all $S \subseteq X$, $\text{rk}(\text{Sp}(S)) = \text{rk}(S)$.
- M12. *Idempotence of Spanning Operator.* For all $S \subseteq X$, $\text{Sp}(\text{Sp}(S)) = \text{Sp}(S)$. (This is theorem T4.)
- M13. *Transitivity of Spanning Operator.* For all $S, T \subseteq X$, if $S \subseteq \text{Sp}(T)$ then $\text{Sp}(S) \subseteq \text{Sp}(T)$. (This is axiom A2.)

It can be shown that the spanning operator of any matroid \mathcal{I} on X satisfies axioms A0, A1, A2, and A3 of §16.1 and that the independent sets with respect to this spanning operator (defined by D4) coincide with the subsets in \mathcal{I} . Conversely, given a spanning operator on a finite set X satisfying A0 through A3, the collection \mathcal{I} of independent sets (defined by D4) can be shown to be a matroid whose spanning operator (defined by (16.1)) coincides with the given spanning operator.

16.15 Equivalence of Matroid Axioms

To illustrate some arguments using definitions H1 through H6 and axioms M1 through M13, we will prove the equivalence of M1, M2, M5, M8, and M10 for a hereditary system \mathcal{I} on a finite set X .

Proof of M1 \Rightarrow M2: Assume M1 holds. Fix $S \subseteq X$, and let T and U be maximal independent subsets of S . To get a contradiction, assume $|U| < |T|$. By M1, we can find $z \in T \sim U$ with $U + z$ independent. But $U + z \subseteq S$ and $U \subsetneq U + z$ contradicts maximality of U . So $|U| = |T|$, and M2 holds.

Proof of M2 \Rightarrow M5: Assume M2 holds. Fix $S, T \subseteq X$. Let A be a fixed maximal

independent subset of $S \cap T$ of size a . Extend A to a maximal independent subset B of $S \cup T$ of size b . Applying M2 to the subsets $S \cap T$ and $S \cup T$, we see that $a = \text{rk}(S \cap T)$ and $b = \text{rk}(S \cup T)$. By maximality of A within $S \cap T$, every element of $B \sim A$ must come from $S \sim T$ or $T \sim S$. So, $B \sim A$ is the disjoint union of two sets $B_1 \subseteq S \sim T$ and $B_2 \subseteq T \sim S$. Letting $b_1 = |B_1|$ and $b_2 = |B_2|$, we have $b - a = b_1 + b_2$. Now, $A \cup B_1$ is independent (being a subset of $A \cup B$) and a subset of S , so $a + b_1 = |A \cup B_1| \leq \text{rk}(S)$. Similarly, $A \cup B_2$ is independent and a subset of T , so $a + b_2 = |A \cup B_2| \leq \text{rk}(T)$. We now compute

$$\text{rk}(S \cap T) + \text{rk}(S \cup T) = a + b = a + a + (b - a) = (a + b_1) + (a + b_2) \leq \text{rk}(S) + \text{rk}(T).$$

So M5 holds.

Proof of M5 \Rightarrow M8: Assume M5 holds. To prove M8, let C and D be distinct circuits in X and let $z \in C \cap D$. We need to show that $(C \cup D) - z$ is dependent. To get a contradiction, suppose that this set is independent. We will apply M5 to the sets $S = C$ and $T = D$. Write $c = |C|$, $d = |D|$, $a = |C \cap D|$, and $b = |C \cup D|$. Consultation of a Venn diagram shows that $a + b = c + d$. Since C is a circuit, C is dependent but $C - z$ is independent, so $\text{rk}(C) = c - 1$. For the same reason, $\text{rk}(D) = d - 1$. Since $C \neq D$ and both sets are circuits, we cannot have $C \subseteq D$ or $D \subseteq C$. It follows that $C \cap D$ is a proper subset of C (and of D). So $C \cap D$ is independent, hence $\text{rk}(C \cap D) = |C \cap D| = a$. Finally, $C \cup D$ contains the dependent set C (and D), so is dependent, but we have assumed $(C \cup D) - z$ is independent. Then $\text{rk}(C \cup D) = |C \cup D| - 1 = b - 1$. Now M5 states that

$$a + (b - 1) = \text{rk}(C \cap D) + \text{rk}(C \cup D) \leq \text{rk}(C) + \text{rk}(D) = (c - 1) + (d - 1) = a + b - 2,$$

yielding the contradiction $-1 \leq -2$. So $(C \cup D) - z$ is dependent, and M8 holds.

Proof of M8 \Rightarrow M10: Assume M8 holds. To prove M10, let $S \subseteq X$ be independent and let $z \in X$. If $z \in S$, then $S + z = S$ contains no circuits. Now suppose $z \notin S$. To get a contradiction, assume that $C \neq D$ are two circuits contained in $S + z$. Both circuits must use z , since S is independent. Then M8 says $(C \cup D) - z$ is dependent, which is impossible since this set is a subset of the independent set S .

Proof of M10 \Rightarrow M1: Assume M10 holds. Let S, T be independent sets in X with $|S| < |T|$. We must find $t \in T \sim S$ with $S + t$ independent. Use induction on $n = |T \sim S| \geq 1$. If $|T \sim S| = 1$, then $|T| > |S|$ forces $|S \sim T| = 0$ and $S \subseteq T$. Letting t be the unique element of $T \sim S$, evidently $S + t = T$ is independent. For the induction step, suppose $|T \sim S| = n > 1$ and the result is known for smaller values of $|T \sim S|$. If $S \subseteq T$, we can add any element of $T \sim S$ to S to get a subset of T , which is independent. Otherwise, fix an $s \in S \sim T$ and ask if $T + s$ is independent. If it is, remove any $t \in T \sim S$ from $T + s$ to get another independent set T' with $|T'| = |T|$ and $|T' \sim S| < |T \sim S|$. By induction, there is $t \in T' \sim S \subseteq T \sim S$ with $S + t$ independent. We must still deal with the case where $T + s$ is dependent. By M10, $T + s$ contains a unique circuit C . We must have $s \in C$ (since T is independent), and C must also use an element $t \in T \sim S$ (since S is independent). If $(T + s) - t$ were dependent, it would contain a circuit D , which would be a circuit contained in $T + s$ unequal to C . This is impossible, so $(T + s) - t$ is independent. We can now finish by induction (as before), since $T' = (T + s) - t$ is independent, has the same size as T , and $|T' \sim S| < |T \sim S|$. So M1 holds.

16.16 Summary

Here we review the main elements of the axiomatic theory of independence structures.

1. *Axioms:* An independence structure consists of a set X and a spanning operator Sp . For each $S \subseteq X$, $\text{Sp}(S)$ is a subset of X containing S . If $\text{Sp}(S)$ contains T , then $\text{Sp}(S)$ contains $\text{Sp}(T)$. If x lies in the span of S , then x is in the span of some finite subset S' of S . Finally, if x is in the span of $S + y$ but not in the span of S alone, then y is in the span of $S + x$.
 2. *Definitions:* If $V = \text{Sp}(S)$, then S spans V and V is a subspace of X ; V is finite-dimensional if S is finite. S is dependent if $x \in \text{Sp}(S - x)$ for some $x \in S$; otherwise S is independent. A basis of X is an independent set spanning X . The dimension of X is the cardinality of any basis.
 3. *Main Results:* The spanning operator is idempotent and inclusion-preserving. Subsets of independent sets are independent, while supersets of dependent sets are dependent. Bases for X coincide with maximal independent subsets of X . We can enlarge an independent set by adjoining any element not in the span of that set. Any independent set can be enlarged to a basis of X ; any spanning set for X contains a basis of X ; so bases of X exist. Independent sets are never larger than spanning sets. The cardinality of a basis of X is unique.
 4. *Algebraic Applications:* The axiomatic setup applies when X is a vector space over a field F (or a division ring); in this case, “independence” means linear independence, “spanning” means linear spanning, “basis” means linear basis, and “dimension” means vector-space dimension. The setup also applies when X is a field extension of a field F ; here, “independence” means algebraic independence, “spanning” means transcendent spanning, “basis” means transcendence basis, and “dimension” means transcendence degree over F .
 5. *Hereditary Systems:* A hereditary system on a finite set X is a nonempty collection \mathcal{I} of subsets of X such that $A \in \mathcal{I}$ and $B \subseteq A$ implies $B \in \mathcal{I}$. Hereditary systems are determined by their independent sets (members of \mathcal{I}), their dependent sets (subsets of X not in \mathcal{I}), their bases (maximal independent sets), their circuits (minimal dependent sets), their rank function ($\text{rk}(S)$ is the maximum size of an independent subset of S), and their spanning operator ($\text{Sp}(S)$ consists of S and those $z \in X$ such that $T + z$ is a circuit for some $T \subseteq S$).
 6. *Matroids:* A matroid consists of a finite set X and a hereditary system \mathcal{I} on X satisfying one (hence all) of axioms M1 through M13 in §16.14.
-

16.17 Exercises

1. Show that axiom A4 automatically holds if X is a finite set.
2. Let X be any set, and define $\text{Sp}(S) = \emptyset$ for all $S \subseteq X$. Which axioms for an independence structure hold?
3. Let X be any set, and define $\text{Sp}(S) = S$ for all $S \subseteq X$. Which axioms for an independence structure hold?
4. Let X be any set, and define $\text{Sp}(S) = X$ for all $S \subseteq X$. Which axioms for an independence structure hold?
5. *Uniform Independence Structures.* (a) Let X be any set, and fix $k \in \mathbb{N}^+$. For $S \subseteq X$, define $\text{Sp}(S) = S$ if $|S| < k$ and $\text{Sp}(S) = X$ if $|S| \geq k$. Prove axioms A0

- through A4. (b) Give an example of (X, Sp) satisfying A0 through A3 but not A4.
6. *Independence Structure for a Set Partition.* Suppose I is an index set, $\{T_i : i \in I\}$ are given pairwise disjoint nonempty sets, and $X = \bigcup_{i \in I} T_i$. Given $S \subseteq X$, define $\text{Sp}(S)$ to be the union of all T_j 's such that $S \cap T_j \neq \emptyset$. Prove that axioms A0 through A4 hold.
 7. Let $X = \mathbb{R}$, and for each $S \subseteq X$, let $\text{Sp}(S) = \overline{S}$, the closure of S in the real line. (By definition, given $z \in \mathbb{R}$ and $S \subseteq \mathbb{R}$, $z \in \overline{S}$ iff for all $\epsilon > 0$, the open interval $(z - \epsilon, z + \epsilon)$ intersects S .) Which axioms for an independence structure hold?
 8. Let X be a commutative group under addition. For $S \subseteq X$, let $\text{Sp}(S)$ be the set of all “finite \mathbb{Z} -linear combinations” $c_1 s_1 + \cdots + c_k s_k$ where $k \in \mathbb{N}$, $c_1, \dots, c_k \in \mathbb{Z}$, and $s_1, \dots, s_k \in S$. Which axioms for an independence structure hold?
 9. For (X, Sp) defined in each of the following exercises, describe all subsets $S \subseteq X$ that span X (as defined in D1). (a) Exc. 2; (b) Exc. 3; (c) Exc. 4; (d) Exc. 5; (e) Exc. 6; (f) Exc. 8 with $X = (\mathbb{Z}, +)$.
 10. For (X, Sp) defined in each of the following exercises, describe all subspaces of X (as defined in D2). (a) Exc. 3; (b) Exc. 4; (c) Exc. 5; (d) Exc. 6; (e) Exc. 8 with $X = (\mathbb{Z}, +)$.
 11. For (X, Sp) defined in each of the following exercises, describe all independent subsets of X (as defined in D4). (a) Exc. 2; (b) Exc. 3; (c) Exc. 4; (d) Exc. 5; (e) Exc. 6; (f) Exc. 7.
 12. For (X, Sp) defined in each of the following exercises, describe all bases of X (as defined in D5) and indicate when X is finite-dimensional (as defined in D6). (a) Exc. 2; (b) Exc. 3; (c) Exc. 4; (d) Exc. 5; (e) Exc. 6; (f) Exc. 7.
 13. Let $X = \{1, 2, 3, 4, 5\}$. For each logical property P , find all maximal subsets $S \subseteq X$ satisfying P . In each case, discuss the existence and uniqueness of the maximal subsets, and note whether all maximal subsets have equal size. (a) P is “All elements of S are prime.” (b) P is “For all $z \in S$, $z+1 \notin S$.” (c) P is “There do not exist $x, y \in S$ with $x+y=6$.” (d) P is “ $2 \in S$, and for all $x, y \in S$, $x+y$ is prime.” (e) P is “For all $z \in \mathbb{Z}$, if $z \in S$ then $z+1 \in S$.”
 14. Let $X = \mathbb{Z}$. For each logical property P , find all maximal subsets $S \subseteq X$ satisfying P or explain why none exist. Is the set of maximal subsets empty, finite, countable, or uncountable? (a) P is “ S is finite.” (b) P is “ $S \neq \mathbb{Z}$.” (c) P is “ S is a proper subgroup of \mathbb{Z} .” (d) P is “For all $z \in \mathbb{Z}$, if $z \in S$ then $z+1 \in S$.” (e) P is “For all $z \in S$, $z+1 \notin S$.” (f) P is “For all $x, y \in S$, $xy \geq 0$.”
 15. Prove that if (X, Sp) satisfies A0, T3, and T4, then axiom A2 follows as a theorem.
 16. (a) Construct a spanning operator on the set $X = \{1, 2, 3\}$ where A0, A1, A3, and T3 hold, but A2 fails. (b) For the spanning operator defined in (a), find all maximal independent subsets of X . Do these all have the same size?
 17. Let X be a finite set. (a) Given a spanning operator Sp on X , show that the set \mathcal{I} of independent subsets of X is a hereditary system. Which axioms are invoked in this proof? (b) Conversely, suppose \mathcal{I} is any hereditary system with ground set X . For $S \subseteq X$, define $\text{Sp}(S)$ by (16.1) (where circuits are defined in terms of \mathcal{I} as in §16.13). Prove that three of the four axioms A0, A1, A2, A3 must hold (which three?).
 18. Given an independence structure (X, Sp) , decide (with explanation) whether each statement is true or false. (a) The union of two independent sets must be

- independent. (b) The intersection of two independent sets must be independent.
- (c) The union of two sets spanning X must also span X . (d) Given a chain $\{S_i : i \in I\}$ of spanning sets of X , $\bigcap_{i \in I} S_i$ must also span X . (e) Any two circuits of X must have the same size. (f) For all $y \in X$ and all independent $U \subseteq X$, $y \in \text{Sp}(U)$ iff $U + y$ is dependent.
19. Prove or disprove the following variant of theorem T14: given $S, T \subseteq X$ with S independent and $X = \text{Sp}(T)$, for any $t \in T \sim S$, there exists $s \in S \sim T$ with $(S - s) + t$ independent.
 20. Let (X, Sp) be a finitely generated independence structure with basis B . (a) Prove $C \subseteq X$ is a basis of X iff C is independent and $|C| \geq |B|$. (b) Prove $C \subseteq X$ is a basis iff C spans X and $|C| \leq |B|$. (c) Show that (a) and (b) can be false if X is not finitely generated.
 21. Let (X, Sp) be a finitely generated independence structure. Aided by theorems in the text, prove that the following matroid axioms must hold. In each case, indicate whether your proof requires the hypothesis that X is finite-dimensional.
(a) M1; (b) M2; (c) M3; (d) M4.
 22. Give an example of a hereditary system \mathcal{I} on a finite set X such that all bases of \mathcal{I} have the same size, but axiom M2 is false.
 23. Prove that each structure (Z, \leq) is a poset. (a) Z consists of any set of subsets of a fixed set X , and for $S, T \in Z$, $S \leq T$ means $S \subseteq T$. (b) $Z = \mathbb{N}^+$, and for $a, b \in \mathbb{Z}$, $a \leq b$ means a divides b . (c) $Z = \mathbb{N}^+$, and for $a, b \in \mathbb{Z}$, $a \leq b$ means b divides a . (d) Z is the set of all functions $f : \mathbb{R} \rightarrow \mathbb{R}$, and for $f, g \in Z$, $f \leq g$ means $f(x) \leq g(x)$ for all $x \in \mathbb{R}$. (e) Z is the set of all functions $f : D \rightarrow \mathbb{R}$ for some $D \subseteq \mathbb{R}$, and for f and g in Z with domains D and E (respectively), $f \leq g$ means $D \subseteq E$ and $f(x) = g(x)$ for all $x \in D$.
 24. For each poset in (b) through (e) of Exercise 23, describe all maximal elements of the poset, or explain why none exist.
 25. For each poset in (b) through (e) of Exercise 23, give an example of an infinite chain in that poset, and state whether the chain has an upper bound in the poset.
 26. Let R be a commutative ring, let I be an ideal of R , and let S be a nonempty subset of $R \sim I$. Use Zorn's lemma to prove that the set of ideals J of R with $I \subseteq J$ and $J \cap S = \emptyset$ has a maximal element.
 27. An ideal P in a nonzero commutative ring R is called *prime* iff $P \neq R$ and for all $x, y \in R$, $xy \in P$ implies $x \in P$ or $y \in P$. Use Zorn's lemma to prove that for a given prime ideal P of R , the set of prime ideals Q contained in P has a *minimal* element relative to set inclusion.
 28. (a) Give an example of a commutative group G that does not have any maximal proper subgroups. (b) Suppose we try to use Zorn's lemma to prove that every commutative group has a maximal proper subgroup, by adapting the argument used to show that maximal ideals exist in commutative rings. Exactly where does the argument break down?
 29. Use Zorn's lemma to prove that any poset (Z, \leq) has a maximal chain (this is a chain $C \subseteq Z$ such that any set properly containing C is not a chain).
 30. Use Zorn's lemma to prove this version of the Axiom of Choice: for any set X , let Z be the set of nonempty subsets of X ; then there exists a function (a set of ordered pairs) $f : Z \rightarrow X$ with $f(S) \in S$ for all $S \in Z$.

31. Suppose (X, Sp) is an independence structure, $S \subseteq T \subseteq X$, S is independent, and T spans X . Prove there is a basis B of X with $S \subseteq B \subseteq T$.
32. (a) In Step 1 of the proof of T26, carefully check that (Z, \leq) satisfies the poset axioms. (b) In Step 2 of the proof of T26, carefully check that $f : A \rightarrow C$ is one-to-one and onto.
33. Let (X, Sp) be an independence structure, and let $Y \subseteq X$. For $S \subseteq Y$, define $\text{Sp}_Y(S) = \text{Sp}(S) \cap Y$. (a) Show (Y, Sp_Y) is an independence structure. (b) Show $S \subseteq Y$ is independent relative to Sp iff S is independent relative to Sp_Y . (c) For X finite-dimensional, show $\text{rk}(Y)$ (defined relative to Sp) is the dimension of Y (defined relative to Sp_Y).
34. In the vector space \mathbb{R}^n (regarded as an independence structure), what are the possible sizes of circuits? Give a specific example of a circuit of each of these sizes.
35. Let V be a vector space over a field F , regarded as an independence structure as in §16.9. Use linear algebra definitions to prove matroid axiom M8 holds for V .
36. In the proofs in §16.9, find all steps where we use the hypothesis that every nonzero element of the field F has a multiplicative inverse. Confirm that commutativity of multiplication in F is never used in the proofs.
37. Let X be a field with subfield K and subset S . As in §16.10, define $K(S)$ to be the set of elements of X that can be written as rational expressions in $K \cup S$. (a) Prove that $K(S)$ is a subfield of X with $K \subseteq K(S)$ and $S \subseteq K(S)$. (b) Prove that if L is any subfield of X with $K \subseteq L$ and $S \subseteq L$, then $K(S) \subseteq L$. (c) Deduce that $K(S)$ is the intersection of all subfields of X that contain $K \cup S$.
38. Let X be a field with subfields $K \subseteq E \subseteq L$. (a) Prove every $z \in K$ is algebraic over K . (b) Prove: if $z \in X$ is algebraic over K , then z is algebraic over E . (c) Prove: if E is a finite extension of K , then E is an algebraic extension of K . (d) Prove: if E is an algebraic extension of K and L is an algebraic extension of E , then L is an algebraic extension of K . (Fix $z \in L$, and say z is a root of $\sum_{i=0}^n c_i x^i \in E[x]$. Apply (c) to $K(c_0, c_1, \dots, c_n, z)$.) (e) Prove: The set of all $z \in X$ that are algebraic over K is a subfield of X containing K . (Use (c).)
39. Decide (with explanation) whether each statement is true or false. (a) \mathbb{C} has transcendence degree 2 over \mathbb{R} . (b) \mathbb{Q} is algebraically closed in \mathbb{R} . (c) Every algebraically independent set over a field F must be F -linearly independent. (d) Every linearly independent set over a field F must be algebraically independent over F . (e) If S is algebraically independent over F and L is a subfield of F , then S is algebraically independent over L . (f) If S is algebraically independent over L and L is a subfield of F , then S is algebraically independent over F . (g) The transcendence degree of a field extension $F \subseteq K$ is zero if K is a finite extension of F . (h) The transcendence degree of a field extension $F \subseteq K$ is zero only if K is a finite extension of F . (i) \mathbb{R} has a finite transcendence basis over \mathbb{Q} .
40. Let X be a field with subfield F , and let $S = \{z_1, \dots, z_n\} \subseteq X$ with all z_i distinct. (a) There is an evaluation homomorphism $E : F[x_1, \dots, x_n] \rightarrow X$ such that $E(c) = c$ for $c \in F$ and $E(x_i) = z_i$ for $1 \leq i \leq n$ (see §3.20). Prove S is algebraically independent over F iff $\ker(E) = \{0\}$. (b) Prove S is algebraically independent over F iff the indexed set of monomials $\{z_1^{e_1} \cdots z_n^{e_n} : (e_1, \dots, e_n) \in \mathbb{N}^n\}$ is linearly independent over F .
41. Given a field F , let $F[x_1, \dots, x_n]$ be the polynomial ring in n indeterminates x_1, \dots, x_n with coefficients in F (see §3.20). (a) Prove $\{x_1, \dots, x_n\}$ is algebraically

independent over F . (b) Let K be the set of “formal fractions” g/h with $g, h \in F[x_1, \dots, x_n]$ and $h \neq 0$. With the ordinary definitions of equality, addition, and multiplication of fractions, it can be shown that K is a field. Show $\{x_1, \dots, x_n\}$ is a transcendence basis of K over F .

42. Let $\{u, v\}$ be algebraically independent over a field F . (a) Prove $\{u + v, uv\}$ is algebraically independent over F . (b) Prove $\{u^2 + 2uv + v^2, u^3 + 3u^2v + 3uv^2 + v^3\}$ is algebraically dependent over F . (c) Is $\{u + v^2, u^2 + v\}$ algebraically independent over F ?
43. Given n polynomials $g_1, \dots, g_n \in \mathbb{Q}[x_1, \dots, x_n]$, define $A \in M_n(\mathbb{Q}[x_1, \dots, x_n])$ by setting $A(i, j) = \partial g_j / \partial x_i$ for $1 \leq i, j \leq n$. Prove: If $\det(A) \neq 0$, then g_1, \dots, g_n are algebraically independent over \mathbb{Q} . [Hint: If the g_i 's are algebraically dependent over \mathbb{Q} , choose a nonzero $h \in \mathbb{Q}[y_1, \dots, y_n]$ of minimum total degree in the y_i 's with $h(g_1, \dots, g_n) = 0$. Use the chain rule to compute the partial derivatives of h .]
44. Fix $n \in \mathbb{N}^+$. (a) For $k \geq 1$, let $p_k = x_1^k + x_2^k + \dots + x_n^k \in \mathbb{Q}[x_1, \dots, x_n]$. Use Exercise 43 to show that $\{p_1, p_2, \dots, p_n\}$ is algebraically independent over \mathbb{Q} . (b) Use Exercise 41 to explain why $\{p_1, \dots, p_{n+1}\} \subseteq \mathbb{Q}[x_1, \dots, x_n]$ is algebraically dependent over \mathbb{Q} .
45. In a graph $G = (V, E, \epsilon)$, prove that for all $v, w \in V$, there is a walk from v to w using edges in $E' \subseteq E$ iff there is a path from v to w using edges in E' .
46. For each graph $G = (V, E, \epsilon)$, list all the bases for the independence structure defined in §16.12. (a) $V = \{1, 2, 3\}$, $E = \{a, b, c, d\}$, $\epsilon(a) = \epsilon(b) = \{1, 2\}$, $\epsilon(c) = \epsilon(d) = \{2, 3\}$. (b) $V = \{1, 2, 3, 4\}$, $E = \{a, b, c, d, e\}$, $\epsilon(a) = \{1, 2\}$, $\epsilon(b) = \{1, 3\}$, $\epsilon(c) = \{1, 4\}$, $\epsilon(d) = \{2, 3\}$, $\epsilon(e) = \{2, 4\}$. (c) $V = \{1, 2, 3, 4, 5, 6\}$, $E = \{a, b, c, d, e, f\}$, $\epsilon(a) = \{1, 2\}$, $\epsilon(b) = \{2, 3\}$, $\epsilon(c) = \{1, 3\}$, $\epsilon(d) = \{4, 5\}$, $\epsilon(e) = \{5, 6\}$, $\epsilon(f) = \{4, 6\}$.
47. Let $G = (V, E, \epsilon)$ be a *connected* graph, which means that for all $u, v \in V$, there is a walk from u to v using edges in E . If $|V| = n$, determine the size of any basis for the independence structure constructed from G in §16.12.
48. Let X be a finite set. (a) Show that for each hereditary system \mathcal{I} on X , the collection \mathcal{B} of bases is nonempty and for all $B, C \in \mathcal{B}$, B is not a proper subset of C . (b) Conversely, given any collection \mathcal{B}' of subsets of X satisfying the condition in (a), prove there exists a unique hereditary system \mathcal{I} on X that has \mathcal{B}' as its collection of bases.
49. Let \mathcal{I} be a hereditary system on a finite set X . (a) Prove $\text{rk}(\emptyset) = 0$. (b) Prove: for all $S \subseteq X$ and $z \in X$, $\text{rk}(S) \leq \text{rk}(S + z) \leq \text{rk}(S) + 1$. (c) Prove: for all $S \subseteq T \subseteq X$, $\text{rk}(S) \leq \text{rk}(T)$ and $\text{Sp}(S) \subseteq \text{Sp}(T)$.
50. Let (X, Sp) be an independence structure with X finite. Prove that equation (16.1) is a theorem (where Sp on the left side is the given spanning operator, and circuits on the right side are minimal dependent sets defined via D3).
51. (a) Let A be an $m \times n$ matrix with entries in a field F . Define \mathcal{I} on the set $X = \{1, 2, \dots, n\}$ by saying that a subset $\{i_1, \dots, i_k\}$ of X is in \mathcal{I} iff the ordered list of columns $(A^{[i_1]}, \dots, A^{[i_k]})$ is F -linearly independent. Show that (X, \mathcal{I}) is a matroid (verify any convenient axiom in the list M1 through M13). (b) Given the real-valued matrix

$$A = \begin{bmatrix} 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & -1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 & -1 \end{bmatrix},$$

find all bases and circuits for the matroid (X, \mathcal{I}) constructed in (a).

52. *Dual Matroids.* Let (X, \mathcal{I}) be a matroid. Let \mathcal{J} be the collection of sets S such that $S \subseteq X \sim B$ for some basis B of \mathcal{I} . Prove that (X, \mathcal{J}) is a matroid by verifying any convenient matroid axiom. This is called the *dual matroid* to (X, \mathcal{I}) .
53. Let \mathcal{I} be a hereditary system on the finite set X . Prove matroid axioms M1, M2, M3, and M4 are logically equivalent.
54. Let \mathcal{I} be a hereditary system on the finite set X . Prove $M5 \Rightarrow M6 \Rightarrow M7 \Rightarrow M1$.
55. (a) In the real vector space $X = \mathbb{R}^2$, give an example of subsets $S, T \subseteq X$ where strict inequality holds in M5. (b) Prove that if S, T are subspaces of a finite-dimensional F -vector space X , then equality will hold in M5.
56. Let \mathcal{I} be a hereditary system on the set X satisfying axiom M9. (a) Using (16.1) to define the spanning operator, prove that axioms A0, A1, A2, and A3 hold. (b) Let \mathcal{J} be the set of independent sets for (X, Sp) , defined via D4. Prove $\mathcal{J} = \mathcal{I}$.
57. Let (X, Sp) be an independence structure with X finite. (a) Show that the set \mathcal{I} of independent subsets of X (defined via D4) is a hereditary system satisfying axiom M9. (b) Prove that the spanning operator for \mathcal{I} (defined in (16.1)) coincides with the spanning operator in the given independence structure.

Elements of Module Theory

Introductory treatments of linear algebra often cover the following fundamental concepts involving vector spaces over a field F : axioms for a vector space, subspaces, quotient spaces, direct sums, linear independence, spanning sets, bases, and linear transformations. The goal of this chapter is to cover the same material in a more general context. We replace the field F by a general ring R and consider “vector spaces” over this ring, which are now called *R*-modules. Intuitively, an *R*-module is a set of “vectors” on which we define operations of vector addition and scalar multiplication satisfying certain axioms. In this setting, the “scalars” come from the ring R rather than a field. As we will see, much of the theory of vector spaces extends without change to this more general situation, although the terminology used is a bit different.

However, we warn the reader that certain aspects of the theory of *R*-modules are quite different from what might be expected based on experience with vector spaces. The most glaring example of this phenomenon is that not every *R*-module has a basis; those modules that do have bases are called *free R*-modules. Even if an *R*-module is free, it may possess bases with different cardinalities, something which does not happen for vector spaces over a field. Fortunately, for commutative rings R (and certain other classes of rings), the cardinality of a basis *is* invariant for free *R*-modules.

Before we can even define modules, we must point out another complication that arises when R is non-commutative. Suppose F is a field, V is an F -vector space, $c, d \in F$, and $v \in V$. We have the associative axiom for scalar multiplication, $(cd)v = c(dv)$, which we might also write as $v(cd) = (vc)d$. The point is that, although we usually write scalars to the left of the vectors on which they act, we may equally well write the scalars on the right instead. However, if we replace F by a non-commutative ring R , so that we now have $c, d \in R$, then the two versions of the associative law just written are no longer equivalent. So we must distinguish between the concepts of scalar multiplication on the left and scalar multiplication on the right. This distinction leads to the concepts of *left* and *right R*-modules.

In the rest of this chapter, we define left and right *R*-modules and discuss fundamental constructions for *R*-modules, including submodules, direct products, direct sums, and quotient modules. We also discuss *R*-module homomorphisms, generating sets for a module, and independent sets in a module, which respectively generalize the concepts of linear transformations, spanning sets, and linearly independent sets in vector spaces. We conclude by examining some properties of free *R*-modules.

17.1 Module Axioms

Let R be an arbitrary ring (see §1.2 for the definition of a ring; recall our convention that every ring has a multiplicative identity element, denoted 1_R). We now present the axioms defining a *left R*-module, which are motivated by the corresponding axioms for a vector space

(see Table 1.4). A left R -module consists of a set M , an *addition* operation $+$: $M \times M \rightarrow M$, and a *scalar multiplication* operation $\cdot : R \times M \rightarrow M$, often denoted by juxtaposition. The addition operation on M must satisfy the following axioms:

- (A1) For all $a, b \in M$, $a + b$ is an element of M (closure under addition).
- (A2) For all $a, b, c \in M$, $(a + b) + c = a + (b + c)$ (associativity).
- (A3) For all $a, b \in M$, $a + b = b + a$ (commutativity).
- (A4) There exists an element $0 \in M$ (necessarily unique) such that $a + 0 = a = 0 + a$ for all $a \in M$ (additive identity).
- (A5) For each $a \in M$, there exists an element $-a \in M$ (necessarily unique) such that $a + (-a) = 0 = (-a) + a$ (additive inverses).

In other words, $(M, +)$ is a commutative group (see §1.1). The scalar multiplication operation must satisfy the following axioms:

- (M1) For all $r \in R$ and $m \in M$, $r \cdot m$ is an element of M (closure under scalar multiplication).
- (M2) For all $m \in M$, $1_R \cdot m = m$ (identity law).
- (M3) For all $r, s \in R$ and $m \in M$, $(rs) \cdot m = r \cdot (s \cdot m)$ (left associativity).

Moreover, the addition and multiplication operations are linked by the following distributive laws:

- (D1) For all $r, s \in R$ and $m \in M$, $(r + s) \cdot m = r \cdot m + s \cdot m$ (distributive law for ring addition).
- (D2) For all $r \in R$ and $m, n \in M$, $r \cdot (m + n) = r \cdot m + r \cdot n$ (distributive law for module addition).

We define *subtraction* by setting $a - b = a + (-b)$ for $a, b \in M$. Elements of M may be called *vectors*, and elements of R *scalars*, although this terminology is more often used when discussing vector spaces over a field.

Now we give the corresponding definition of a *right R -module*. A right R -module consists of a set N , an *addition* operation $+: N \times N \rightarrow N$, and a *scalar multiplication* operation $\star : N \times R \rightarrow N$. We require $(N, +)$ to be a commutative group (i.e., axioms (A1) through (A5) must hold with M replaced by N). The multiplication operation must satisfy the following axioms:

- (M1') For all $n \in N$ and $r \in R$, $n \star r$ is an element of N (closure under scalar multiplication).
- (M2') For all $n \in N$, $n \star 1_R = n$ (identity law).
- (M3') For all $n \in N$ and $r, s \in R$, $n \star (rs) = (n \star r) \star s$ (right associativity).

Moreover, the addition and multiplication operations are linked by the following distributive laws:

- (D1') For all $n \in N$ and $r, s \in R$, $n \star (r + s) = n \star r + n \star s$ (distributive law for ring addition).
- (D2') For all $m, n \in N$ and $r \in R$, $(m + n) \star r = m \star r + n \star r$ (distributive law for module addition).

Subtraction is defined as before.

When R is commutative, there is no essential difference between left R -modules and right R -modules. To see why, let $(M, +)$ be a fixed commutative group. We can set up a one-to-one correspondence between “left scalar multiplications” $\cdot : R \times M \rightarrow M$ and “right scalar multiplications” $\star : M \times R \rightarrow M$, given by $m \star r = r \cdot m$ for all $r \in R$ and $m \in M$. For any ring R , one may verify that the left module axioms (M1), (M2), (D1), and (D2) hold for \cdot iff the corresponding right module axioms (M1'), (M2'), (D1'), and (D2') hold for \star . If R is commutative, then (M3) holds for \cdot iff (M3') holds for \star . Thus, in the commutative case, it makes little difference whether we multiply by scalars on the left or on the right. Later in this chapter, the unqualified term “module” will always mean *left* R -module; but there will be analogous definitions and results for right R -modules.

Next we introduce the concept of *R -module homomorphisms*, which are analogous to linear transformations of vector spaces. Let M and N be left R -modules. A map $f : M \rightarrow N$ is a *left R -module homomorphism* iff $f(x+y) = f(x) + f(y)$ and $f(rx) = rf(x)$ for all $x, y \in M$ and $r \in R$. The definition of a *right R -module homomorphism* between right R -modules M and N is analogous: we require $f(x+y) = f(x) + f(y)$ and $f(xr) = f(x)r$ for all $x, y \in M$ and $r \in R$. Homomorphisms of R -modules (left or right) are also called *R -maps* or *R -linear maps*. The following terminology is sometimes used for R -maps satisfying additional conditions: a module homomorphism f is called a *monomorphism*, *epimorphism*, *isomorphism*, *endomorphism*, or *automorphism*, iff f is injective, f is surjective, f is bijective, $M = N$, or f is bijective and $M = N$ (respectively). The modules M and N are called *isomorphic* or *R -isomorphic* iff there exists an R -module isomorphism $g : M \rightarrow N$; in this case, we write $M \cong N$. We let the reader verify that the following facts about linear transformations and group homomorphisms are also true for R -maps: the identity map on M is a bijective R -map; the composition of R -maps is an R -map (similarly for monomorphisms, epimorphisms, etc.); the inverse of an R -isomorphism is an R -isomorphism; the relation “ M is isomorphic to N as an R -module” is an equivalence relation on any collection of left R -modules; and for an R -map $f : M \rightarrow N$, $f(0_M) = 0_N$ and $f(-x) = -f(x)$ for all $x \in M$. One may also check that $0_R x = 0_M$ for all $x \in M$; $r 0_M = 0_M$ for all $r \in R$; and $r(-x) = (-r)x = -(rx)$ for all $r \in R$ and all $x \in M$.

17.2 Examples of Modules

We now discuss four fundamental examples of R -modules.

First, any ring R is a left R -module, if we take addition and scalar multiplication to be the given addition and multiplication in the ring R . In this case, the module axioms reduce to the axioms in the definition of a ring (see §1.2). Similarly, every ring R is a right R -module. Suppose $f : R \rightarrow R$ is a function satisfying $f(x+y) = f(x) + f(y)$ for all $x, y \in R$. Then f is a ring homomorphism iff $f(1_R) = 1_R$ and $f(xy) = f(x)f(y)$ for all $x, y \in R$. On the other hand, f is a left R -module homomorphism iff $f(xy) = xf(y)$ for all $x, y \in R$, while f is a right R -module homomorphism iff $f(xy) = f(x)y$ for all $x, y \in R$.

Second, generalizing the first example, let R be a subring of a ring S . (Recall from §1.4 that this means R is a subset of S containing 0_S and 1_S and closed under addition, subtraction, and multiplication; R itself is a ring under the operations inherited from S .) The ring S is a left R -module, if we take addition to be the addition in S and scalar multiplication $\cdot : R \times S \rightarrow S$ to be the restriction of the multiplication $\cdot : S \times S \rightarrow S$ in the ring S . Here, the conditions in the module axioms hold because they are a subset of the conditions in the ring axioms for S . Similarly, S is a right R -module.

Third, let F be a field. Comparing the module axioms to the axioms defining a vector space over F (see Table 1.4), we see that a left F -module V is exactly the same as an F -vector space. An F -module homomorphism $T : V \rightarrow W$ is exactly the same as an F -linear map from V to W .

Fourth, let $(M, +)$ be any commutative group, and consider the ring $R = \mathbb{Z}$. We make M into a left \mathbb{Z} -module by defining $0 \cdot x = 0_M$, $n \cdot x = x + x + \cdots + x$ (n copies of x), and $-n \cdot x = -(x + x + \cdots + x)$ (n copies of x) for all $n > 0$ and $x \in M$. The module axioms are routinely checked; in fact, the axioms regarding multiplication are additive versions of the “laws of exponents” for commutative groups. Moreover, one can verify (using axioms (M2) and (D1)) that \cdot is the *unique* scalar multiplication map from $\mathbb{Z} \times M$ into M that makes $(M, +)$ into a left \mathbb{Z} -module. In other words, the “ \mathbb{Z} -module structure” of M is completely determined by M ’s structure as an additive group. Moreover, consider a function $f : M \rightarrow N$ between two \mathbb{Z} -modules. On one hand, if f is a \mathbb{Z} -module homomorphism, then it is in particular a homomorphism of commutative groups (meaning $f(x + y) = f(x) + f(y)$ for all $x, y \in M$). Conversely, if f is a group homomorphism, then we know from group theory that $f(n \cdot x) = n \cdot f(x)$ for all $x \in M$ and $n \in \mathbb{Z}$, so that f is automatically a \mathbb{Z} -module homomorphism. These remarks show that \mathbb{Z} -modules and \mathbb{Z} -module homomorphisms are essentially the same as commutative groups and group homomorphisms.

Here we begin to see one of the advantages of module theory as a unifying notational tool: facts about vector spaces, commutative groups, and rings can all be discussed in the general framework of module theory.

17.3 Submodules

The next few sections discuss some general constructions for manufacturing new modules from old ones: submodules, direct products, direct sums, Hom modules, quotient modules, and changing the ring of scalars. Chapter 20 discusses more advanced constructions for modules, such as tensor products, tensor powers, exterior powers, and symmetric powers.

Let N be a subset of a left R -module M . We say N is an *R -submodule* of M (or simply a *submodule* if R is understood) iff N is a subgroup of the additive group $(M, +)$ such that for all $r \in R$ and $x \in N$, $r \cdot x \in N$. In other words, a submodule of M is a subset of M containing 0_M and closed under addition, additive inverses, and left multiplication by scalars from R . Using the fact that $-1 \in R$, one sees that closure under inverses follows from the other closure conditions, since $-n = -(1_R \cdot n) = (-1_R) \cdot n$ for $n \in N$. If we restrict the addition and scalar multiplication operations for M to the domains $N \times N$ and $R \times N$ (respectively), then N becomes a left R -module. The axioms (A1) and (M1) hold for N by definition of a submodule. The other axioms for N are special cases of the corresponding axioms for M , keeping in mind that $0_M \in N$ and $-x \in N$ whenever $x \in N$.

The definition of a submodule N of a right R -module M is similar: we now require N to be an additive subgroup of M such that $x \star r \in N$ for all $x \in N$ and $r \in R$.

Consider the examples in the preceding section. Regarding R as a left R -module, the definition of a left R -submodule of R is exactly the definition of a *left ideal* of R (see §1.4). Regarding R as a right R -module, the definition of a right R -submodule of R is exactly the definition of a *right ideal* of R . If R is commutative, the R -submodules of R (viewed as a left or right R -module) are precisely the *ideals* of the ring R . Next, if V is a vector space over a field F (so V is an F -module), the F -submodules of V are precisely the *subspaces* of the vector space V . Finally, if M is a commutative group and hence a \mathbb{Z} -module, the \mathbb{Z} -submodules of M are precisely the *subgroups* of M . For, if H is a subgroup, $x \in H$, and

$n \in \mathbb{Z}$, then $n \cdot x \in H$ follows by induction from the fact that H is closed under addition and additive inverses.

We now discuss intersections and sums of submodules. Let $\{M_i : i \in I\}$ be a nonempty family of submodules of an R -module M . The *intersection* $N = \bigcap_{i \in I} M_i$ is a submodule of M , hence an R -module. For, $0_M \in N$ since 0_M lies in every M_i . If $x, y \in N$, then $x, y \in M_i$ for all $i \in I$, so $x + y \in M_i$ for all $i \in I$, so $x + y \in N$. If $x \in N$ and $r \in R$, then $x \in M_i$ for all $i \in I$, so $r \cdot x \in M_i$ for all $i \in I$, so $r \cdot x \in N$. N is the largest submodule of M contained in every M_i .

Next let M and N be submodules of an R -module P . One may verify that the set $M + N = \{m + n : m \in M, n \in N\}$ is a submodule of P , called the *sum* of M and N . Sums of finitely many submodules of P are defined analogously. For the general case, let $\{M_i : i \in I\}$ be a family of submodules of P . The *sum* of these submodules, denoted $\sum_{i \in I} M_i$, is the set of all *finite* sums $m_{i_1} + \cdots + m_{i_s}$ where $i_j \in I$, $s \in \mathbb{N}$, and $m_{i_j} \in M_{i_j}$. One may check that $\sum_{i \in I} M_i$ is an R -submodule of P , and this is the smallest submodule of P containing every M_i . If I is empty, we consider the vacuous sum $\sum_{i \in I} M_i$ to be the submodule $\{0_M\}$ consisting of zero alone.

A *partially ordered set* is a set X and a relation \leq on X that is *reflexive* (for all $x \in X$, $x \leq x$), *antisymmetric* (for all $x, y \in X$, if $x \leq y$ and $y \leq x$ then $x = y$), and *transitive* (for all $x, y, z \in X$, if $x \leq y$ and $y \leq z$ then $x \leq z$). A *lattice* is a partially ordered set in which any two elements have a least upper bound and a greatest lower bound (see the Appendix for more detailed definitions). A *complete lattice* is a partially ordered set (X, \leq) in which any nonempty subset of X has a least upper bound and a greatest lower bound. In more detail, given any nonempty $S \subseteq X$, the set of upper bounds $\{x \in X : s \leq x \text{ for all } s \in S\}$ must be nonempty and have a least element, and similarly for the set of lower bounds of S . For any left R -module M , the set X of all submodules of M is a partially ordered set under the ordering defined by $N \leq P$ iff $N \subseteq P$ (for $N, P \in X$). Our results above show that the poset X of submodules of a left R -module is a complete lattice, since every nonempty family of submodules $\{M_i : i \in I\}$ has a greatest lower bound (the intersection of the M_i) and a least upper bound (the sum of the M_i).

A left R -module M is called *simple* iff $M \neq \{0\}$ and the only R -submodules of M are $\{0\}$ and M . The adjective “simple” really describes the submodule lattice of M , which is a poset with only two elements. For example, \mathbb{Z}_2 , \mathbb{Z}_3 , and \mathbb{Z}_5 are simple \mathbb{Z} -modules, but \mathbb{Z}_1 and \mathbb{Z}_4 are not simple \mathbb{Z} -modules. One can check that a commutative ring R is simple (as a left R -module) iff R is a field.

17.4 Submodule Generated by a Subset

Let S be any subset of a left R -module M , not necessarily a submodule. We would like to construct a submodule of M containing S that is as small as possible. To do this, let $\{M_i : i \in I\}$ be the family of all submodules of M that contain S . This family is nonempty, since M itself is a submodule of M containing S . Let $N = \bigcap_{i \in I} M_i$. Then N is a submodule, N contains S , and N is contained in any submodule of M that contains S (which must be one of the M_i 's). We write $N = \langle S \rangle$, and call N the *submodule generated by S* . Elements of S are called *generators* for N . In the case of vector spaces, we would say that the subspace N is *spanned* by the vectors in S , which form a *spanning set* for N .

Now we give another description of $\langle S \rangle$ that shows how elements of this submodule are “built up” from elements of S . Let N' be the set of all finite sums $r_1s_1 + \cdots + r_ns_n$, where $r_i \in R$, $s_i \in S$, and $n \in \mathbb{N}$. Such a sum is called an *R -linear combination* of elements of N .

S . (If $n = 0$, the sum is defined to be zero, so 0_M is always in N' .) We will show that $N' = N = \langle S \rangle$. First, one can check that N' is an R -submodule of M . Since R has an identity, each $s \in S$ can be written as $1_R \cdot s$, which implies that N' contains S . Therefore, N' is one of the submodules M_i appearing in the definition of N , and hence $N' \supseteq N$. Conversely, consider an arbitrary submodule M_i that contains S . Since M_i is a submodule, it must contain $r_1 s_1 + \cdots + r_n s_n$, which is the typical element of N' . Thus, $N' \subseteq M_i$ for every M_i , and hence $N' \subseteq N$. Some special cases of this result should be mentioned separately. If $S = \{s_1, \dots, s_k\}$ is finite, then $\langle S \rangle = \{\sum_{i=1}^k r_i s_i : r_i \in R\}$. If $S = \{s\}$ has one element, then $\langle S \rangle = \{rs : r \in R\}$. If S is empty, then $\langle S \rangle = \{0_M\}$.

Given any R -module M , we say M is *finitely generated* iff there exists a finite subset S of M such that $M = \langle S \rangle$. We say M is *cyclic* iff there exists a one-element set S such that $M = \langle S \rangle$. If $S = \{x\}$, we often write $M = \langle x \rangle$ or $M = Rx$ instead of $M = \langle \{x\} \rangle$.

One property of generating sets is that *module homomorphisms that agree on a generating set must be equal*. More formally, let M be an R -module generated by S , and let $f, g : M \rightarrow P$ be R -module homomorphisms such that $f(s) = g(s)$ for all $s \in S$. Then $f = g$. To prove this, take any nonzero element x of M and write it as $x = r_1 s_1 + \cdots + r_k s_k$ with $r_i \in R$ and $s_i \in S$. Since f and g are module homomorphisms agreeing on S , $f(x) = \sum_{i=1}^k r_i f(s_i) = \sum_{i=1}^k r_i g(s_i) = g(x)$. Also, $f(0_M) = 0_P = g(0_M)$, so that $f = g$.

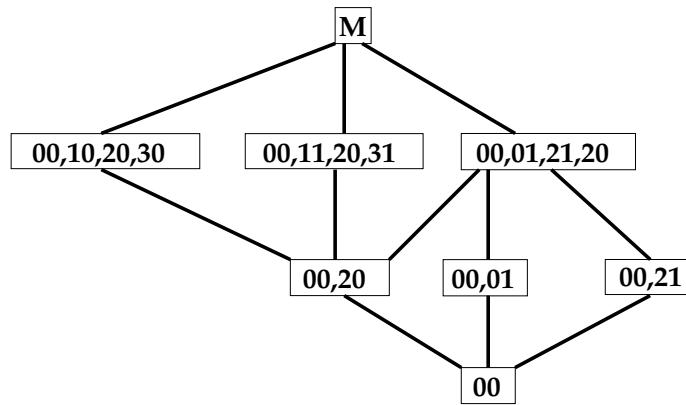
Let us use the ideas of cyclic submodules and generating sets to determine all the \mathbb{Z} -submodules of the \mathbb{Z} -module $M = \mathbb{Z}_4 \times \mathbb{Z}_2$. M is an 8-element commutative group, and its submodules are precisely the subgroups of M . It is known from group theory (Lagrange's theorem) that the size of each such subgroup must be a divisor of 8, namely 1 or 2 or 4 or 8. We can start finding subgroups by looking at the cyclic submodules generated by each element of M . We discover the following submodules:

$$\begin{aligned} A &= \mathbb{Z}(0, 0) = \{(0, 0)\} \\ B &= \mathbb{Z}(0, 1) = \{(0, 0), (0, 1)\} \\ C &= \mathbb{Z}(1, 0) = \{(0, 0), (1, 0), (2, 0), (3, 0)\} = \mathbb{Z}(3, 0) \\ D &= \mathbb{Z}(2, 0) = \{(0, 0), (2, 0)\} \\ E &= \mathbb{Z}(2, 1) = \{(0, 0), (2, 1)\} \\ F &= \mathbb{Z}(1, 1) = \{(0, 0), (1, 1), (2, 0), (3, 1)\} = \mathbb{Z}(3, 1). \end{aligned}$$

We can build further submodules by taking sums of the cyclic submodules found above. This produces two new submodules, namely $G = B + D = \{(0, 0), (0, 1), (2, 0), (2, 1)\}$ and $B + C = \mathbb{Z}_4 \times \mathbb{Z}_2 = M$. Figure 17.1 displays the lattice of submodules of the \mathbb{Z} -module M . In the figure, we abbreviate (x, y) as xy , and we draw a thick line from a submodule U up to a submodule V whenever $U \subseteq V$ and there is no submodule properly between U and V .

17.5 Direct Products, Direct Sums, and Hom Modules

Let I be any set, and suppose we have a left R -module M_i for each $i \in I$. Let N be the set of all functions $f : I \rightarrow \bigcup_{i \in I} M_i$ such that $f(i) \in M_i$ for all $i \in I$. It may be helpful to think of such a function as an “ I -tuple” $(f(i) : i \in I)$, particularly when I is a finite set such as $\{1, 2, \dots, n\}$. Given $f, g \in N$, define their *sum* $f + g$ by the rule $(f + g)(i) = f(i) + g(i) \in M_i$ for $i \in I$. One can check that N becomes a commutative group with this addition operation. Next, for $f \in N$ and $r \in R$, define $rf \in N$ by the rule $(rf)(i) = r \cdot [f(i)] \in M_i$ for $i \in I$. One can verify that this scalar multiplication makes N into a left R -module. The axioms

**FIGURE 17.1**Lattice of Submodules of the \mathbb{Z} -Module $\mathbb{Z}_4 \times \mathbb{Z}_2$.

for N follow from the corresponding axioms for the M_i 's. For instance, to check (D1), fix $r, s \in R$ and $f \in N$, and compare the effects of the functions $(r + s)f$ and $rf + sf$ on a typical $i \in I$:

$$\begin{aligned} [(r + s)f](i) &= (r + s) \cdot [f(i)] = r \cdot [f(i)] + s \cdot [f(i)] \text{ (by (D1) in } M_i) \\ &= (rf)(i) + (sf)(i) = (rf + sf)(i). \end{aligned}$$

We write $N = \prod_{i \in I} M_i$ and call N the *direct product* of the modules M_i . If $I = \{1, 2, \dots, n\}$, we write $N = M_1 \times M_2 \times \dots \times M_n$. In this case, elements of N are usually thought of as “ n -tuples” instead of functions. Using the n -tuple notation for elements of N , the module operations are

$$(f_1, \dots, f_n) + (g_1, \dots, g_n) = (f_1 + g_1, \dots, f_n + g_n) \quad (f_i, g_i \in M_i);$$

$$r(f_1, \dots, f_n) = (rf_1, \dots, rf_n) \quad (f_i \in M_i, r \in R),$$

where we have written f_i instead of $f(i)$.

As a special case of the direct product, we can take every M_i equal to a given module M , and we may write $N = M^I$ in this case. If $I = \{1, 2, \dots, n\}$, we often write $N = M^n$. For instance, the vector spaces F^n (for F a field) are special cases of this construction. Taking $R = \mathbb{Z}$, we see that direct products of commutative groups are also special cases. There is a similar construction for producing the direct product of an indexed set of right R -modules.

Consider again the direct product $N = \prod_{i \in I} M_i$ of left R -modules. We say that a function $f \in N$ has *finite support* iff the set $\{i \in I : f(i) \neq 0_{M_i}\}$ is finite. Let N_0 consist of all functions $f \in N$ with finite support. One may check that N_0 contains 0_N and is closed under addition and scalar multiplication, so that N_0 is an R -submodule of N . We call N_0 the *direct sum* of the R -modules M_i , denoted $\bigoplus_{i \in I} M_i$. If I is finite, the direct sum coincides with the direct product; so, for example, $M_1 \times M_2 \times \dots \times M_n = M_1 \oplus M_2 \oplus \dots \oplus M_n$.

Let F be a field. We show that *every F-vector space V is isomorphic to a direct sum of copies of F*. The proof requires the well-known fact (proved in §16.9) that every vector space has a basis. Given V , let X be a fixed basis of V , and let W be the direct sum of copies of F indexed by X . Thus, a typical element of W is a function $g : X \rightarrow F$ with finite support. To get a vector space isomorphism $T : W \rightarrow V$, define $T(g) = \sum_{x \in X} g(x)x \in V$ for $g \in W$. The sum makes sense because, for each fixed $g \in W$, there are only finitely

many nonzero summands. To describe the inverse map $S : V \rightarrow W$, recall that each $v \in V$ can be uniquely written as a finite linear combination $v = c_1x_1 + \cdots + c_nx_n$ for some $n \geq 0$, nonzero $c_i \in F$, and distinct $x_i \in X$. Define $S(v)$ to be the function $g : X \rightarrow F$ such that $g(x_i) = c_i$ for $1 \leq i \leq n$, and $g(x) = 0$ for all other $x \in X$. One verifies that S and T are linear maps that are inverses of each other. We will see later that this result does not extend to general R -modules, because not every R -module has a basis.

Now let M and N be left R -modules. We write $\text{Hom}(M, N)$ to denote the set of all group homomorphisms $f : M \rightarrow N$. Note that $\text{Hom}(M, N)$ is a subset of the direct product $N^M = \prod_{x \in M} N$, which is a left R -module. We check that this subset is, in fact, an R -submodule of N^M . The zero element of N^M is the zero map $z : M \rightarrow N$ such that $z(x) = 0_N$ for all $x \in M$. This map is evidently a group homomorphism, so $0_{N^M} \in \text{Hom}(M, N)$. Next, fix $f, g \in \text{Hom}(M, N)$. Then $f + g \in \text{Hom}(M, N)$ because for $x, y \in M$,

$$\begin{aligned} (f + g)(x + y) &= f(x + y) + g(x + y) = f(x) + f(y) + g(x) + g(y) \\ &= f(x) + g(x) + f(y) + g(y) \quad (\text{since } + \text{ in } N \text{ is commutative}) \\ &= (f + g)(x) + (f + g)(y). \end{aligned}$$

Suppose $f \in \text{Hom}(M, N)$ and $r \in R$. We have $rf \in \text{Hom}(M, N)$ because for $x, y \in M$,

$$(rf)(x + y) = r \cdot [f(x + y)] = r \cdot [f(x) + f(y)] = r \cdot f(x) + r \cdot f(y) = (rf)(x) + (rf)(y).$$

Since $\text{Hom}(M, N)$ is an R -submodule of N^M , it is in particular a left R -module under the pointwise operations on functions inherited from N^M .

Now consider the set $\text{Hom}_R(M, N)$ of all left R -module homomorphisms $g : M \rightarrow N$, which is another subset of N^M (and also a subset of $\text{Hom}(M, N)$). We show that, if R is commutative, then $\text{Hom}_R(M, N)$ is a submodule of N^M and, hence, is itself a left R -module. One checks, as above, that $\text{Hom}_R(M, N)$ is an additive subgroup of N^M (for any ring R). Fix $r \in R$, $f \in \text{Hom}_R(M, N)$, $s \in R$, and $x \in M$. For R commutative, we have

$$\begin{aligned} (rf)(s \cdot x) &= r \cdot [f(s \cdot x)] = r \cdot [s \cdot f(x)] = (rs) \cdot [f(x)] \\ &= (sr) \cdot [f(x)] \quad (\text{by commutativity of } R) \\ &= s \cdot (r \cdot [f(x)]) = s \cdot ((rf)(x)). \end{aligned}$$

The other condition for being an R -map, namely $(rf)(x+y) = (rf)(x) + (rf)(y)$ for $x, y \in M$, already follows from our calculations in the last paragraph. Therefore, $rf \in \text{Hom}_R(M, N)$, as needed.

As a special case of this construction, consider F -vector spaces V and W . The set of all linear transformations from V to W , sometimes denoted by $L(V, W)$, becomes an F -vector space under pointwise operations on functions. This follows by noting that $L(V, W) = \text{Hom}_F(V, W)$ and F is commutative.

17.6 Quotient Modules

In §1.6, we discussed the construction of the quotient group of a commutative group by a subgroup, and the formation of the quotient space of a vector space by a subspace. These constructions are special cases of quotient modules, which are defined as follows.

Let N be a submodule of a left R -module M . For each $x \in M$, we have the *coset* $x+N = \{x+n : n \in N\}$. The *quotient set* is the collection of cosets $M/N = \{x+N : x \in M\}$.

We saw in §1.6 that for all $x, z \in M$, $x+N = z+N$ iff $x-z \in N$. Moreover, M is the disjoint union of the distinct cosets of N , i.e., $M = \bigcup_{S \in M/N} S$. Since N is a subgroup of $(M, +)$, we can define an addition operation on M/N by setting $(x+N) + (y+N) = (x+y)+N$ for all $x, y \in M$. We proved in §1.6 that this binary operation is well-defined and satisfies the axioms for a commutative group. The identity element of this quotient group is $0_{M/N} = 0_M + N = \{0+n : n \in N\} = N$, and the additive inverse of the coset $x+N$ is the coset $(-x)+N$, for any $x \in M$.

To complete the definition of the quotient R -module M/N , we introduce a scalar multiplication operation $\cdot : R \times M/N \rightarrow M/N$. Given $r \in R$ and $x \in M$, define $r \cdot (x+N) = (rx)+N$. We must check that this operation is well-defined. Say $r \in R$ and $x, z \in M$ are such that $x+N = z+N$; is it true that $r \cdot (x+N) = r \cdot (z+N)$? First of all, $x-z \in N$. Since N is an R -submodule, $r(x-z) \in N$. So $rx - rz \in N$, and hence $(rx)+N = (rz)+N$. This shows that $r \cdot (x+N) = r \cdot (z+N)$. Now that we know scalar multiplication is well-defined, it is routine to check that M/N is a left R -module. For example, the following calculation verifies axiom (D1): for $x \in M$ and $r, s \in R$,

$$\begin{aligned} (r+s) \cdot (x+N) &= ((r+s)x+N) \\ &= ((rx+sx)+N) \quad (\text{by (D1) in } M) \\ &= (rx+N) + (sx+N) = [r \cdot (x+N)] + [s \cdot (x+N)]. \end{aligned}$$

For any left R -module M and submodule N , there is a surjection $p : M \rightarrow M/N$ defined by $p(x) = x+N$ for $x \in M$. This map is called the *natural map* from M to M/N , the *canonical map* from M to M/N , or the *projection* of M onto M/N . The map p is R -linear. For, by definition of the operations in M/N ,

$$p(x+y) = (x+y)+N = (x+N) + (y+N) = p(x) + p(y), \text{ and}$$

$$p(rx) = (rx)+N = r \cdot (x+N) = r \cdot p(x)$$

for all $x, y \in M$ and $r \in R$.

Let us consider some special cases and extensions of the quotient module construction. For $R = \mathbb{Z}$, the quotient \mathbb{Z} -module M/N is really the same thing as the quotient group M/N (since, as we have seen, the \mathbb{Z} -module structure on M/N is uniquely determined by the addition operation on this set). If R is a field F , then N is a vector subspace of the F -vector space M , and M/N is precisely the quotient vector space of M by N . Finally, suppose R is arbitrary, and M is R viewed as a left R -module. Given any left ideal I of R , R/I is a left R -module. However, more can be said if I is a (two-sided) ideal of R (i.e., if I is both a left ideal and a right ideal, which always happens for commutative rings R). By analogy with the above constructions, we can define a multiplication operation $\cdot : R/I \times R/I \rightarrow R/I$ (which is distinct from the scalar multiplication defined above) by setting $(a+I) \cdot (b+I) = (ab)+I$ for all $a, b \in R$. This operation is well-defined, for if $a+I = a'+I$ and $b+I = b'+I$ (where $a, b, a', b' \in R$), then $(ab)+I = (a'b')+I$ because

$$a-a' \in I, \quad b-b' \in I, \quad ab - a'b' = a(b-b') + (a-a')b',$$

and the latter expression lies in I because I is a two-sided ideal. With this multiplication and the addition already defined, one checks that R/I becomes a ring, which is commutative if R is commutative. R/I is the *quotient ring of R by the ideal I* .

One of the nice features of modules is that the quotient of a module by any submodule is always defined, unlike the situation for general groups (where the subgroup must be normal) or rings (where the subset must be a two-sided ideal).

Suppose S is a generating set for the R -module M , and N is a submodule of M . Then the

image of S in M/N , namely $p[S] = \{s + N : s \in S\}$, is a generating set for M/N . To prove this, note that any coset in M/N has the form $x + N$ for some $x \in M$. Write $x = \sum_i a_i s_i$ with $a_i \in R$ and $s_i \in S$. Then $x + N = (\sum_i a_i s_i) + N = \sum_i ((a_i s_i) + N) = \sum_i a_i (s_i + N)$.

17.7 Changing the Ring of Scalars

This section describes two constructions for changing the ring of scalars for a given module. First, let M be a left R -module, and suppose S is a subring of R . Restricting the scalar multiplication $\cdot : R \times M \rightarrow M$ to the domain $S \times M$, we get a scalar multiplication of S on M that satisfies the module axioms. Hence, we can regard M as a left S -module. More generally, if S is any ring and $f : S \rightarrow R$ is a ring homomorphism, then one can verify that the scalar multiplication $\star : S \times M \rightarrow M$, defined by $s \star x = f(s) \cdot x$ for $s \in S$ and $x \in M$, turns M into a left S -module. The situation where S is a subring is the special case where $f : S \rightarrow R$ is the inclusion map given by $f(s) = s$ for all $s \in S$.

Second, let M be a left R -module. Suppose I is an ideal of R , and let S be the quotient ring R/I . If I has the property that $i \cdot x = 0_M$ for all $i \in I$ and $x \in M$, then we can regard M as an S -module (i.e., an R/I -module) by defining $(r+I) \bullet x = r \cdot x$ for $r \in R$ and $x \in M$. The italicized condition, which we express by saying that I annihilates M , is needed to show that this new scalar multiplication operation is well-defined. For, fix $r, r' \in R$ and $x \in M$ with $r+I = r'+I$. Then $r - r' \in I$ implies that $(r - r') \cdot x = 0$ (by the annihilation condition), so that $r \cdot x = r' \cdot x$ and therefore $(r+I) \bullet x = (r'+I) \bullet x$. Once we know that the operation is well-defined, the S -module axioms follow routinely from the corresponding R -module axioms. For instance, (M2) follows since $(1_R + I) \bullet x = 1_R \cdot x = x$ for all $x \in M$.

Now let N be any R -module, not necessarily annihilated by the ideal I . Define a subset IN of N , which consists of all finite sums of terms $i \cdot x$ with $i \in I$ and $x \in N$. One can check that IN is an R -submodule of N , so that we can form the quotient module $M = N/IN$. Now, M is annihilated by I , because $i \cdot (x + IN) = (i \cdot x) + IN = 0_N + IN = 0_M$ for all $i \in I$ and $x \in N$ (since $i \cdot x - 0_N \in IN$). Therefore, M is an S -module. We restate this result for emphasis: given any R -module N and any ideal I of R , N/IN is an R/I -module via the rule $(r+I) \bullet (x + IN) = (r \cdot x) + IN$ for $r \in R$ and $x \in N$.

Here is one situation in which this result can be helpful. Let R be a nonzero commutative ring, N an R -module, and I a maximal ideal of R (i.e., there are no ideals J with $I \subsetneq J \subsetneq R$). We assume the well-known facts that such maximal ideals exist in R , and maximality of I implies $F = R/I$ is a field (for proofs, see §16.6 and Exercise 36). Hence, we can pass from the general R -module N to an associated *vector space* over F , namely, N/IN . As we will see later, this association lets us use known theorems about vector spaces over a field to deduce facts about modules over a commutative ring.

17.8 Fundamental Homomorphism Theorem for Modules

The next two sections describe some central results concerning R -module homomorphisms, which parallel the corresponding theory for homomorphisms of groups, rings, and vector spaces (cf. §1.7).

Let $f : M \rightarrow N$ be a homomorphism of left R -modules. The *kernel* of f , denoted $\ker(f)$, is the set $\{x \in M : f(x) = 0_N\} \subseteq M$. The *image* of f , denoted $\text{img}(f)$, is the

set $\{y \in N : y = f(x) \text{ for some } x \in M\} \subseteq N$. One can check that $\ker(f)$ is a submodule of M , and $\text{img}(f)$ is a submodule of N . More generally, if M' is any submodule of M , then the direct image $f[M'] = \{f(x) : x \in M'\}$ is a submodule of N . The image of f is the special case where $M' = M$. If N' is any submodule of N , then the inverse image $f^{-1}[N'] = \{x \in M : f(x) \in N'\}$ is a submodule of M . The kernel of f is the special case obtained by taking $N' = \{0_N\}$. Note that $f^{-1}[N']$ contains $\ker(f)$ for any submodule N' of N .

From the definition, we see that f is surjective iff $\text{img}(f) = N$. Moreover, we claim that f is injective iff $\ker(f) = \{0_M\}$. For, suppose the kernel is zero. Assume $x, z \in M$ satisfy $f(x) = f(z)$. Then $f(x - z) = f(x) - f(z) = 0$, so that $x - z \in \ker(f) = \{0\}$. Therefore, $x = z$, so that f is injective. Conversely, assume f is injective. On one hand, $0_M \in \ker(f)$ since $f(0_M) = 0_N$. On the other hand, for any $x \in M$ unequal to 0, injectivity of f gives $f(x) \neq f(0) = 0$, so $x \notin \ker(f)$.

Using kernels, images, and quotient modules, we can build an R -module *isomorphism* from any given R -module *homomorphism*. This is the purpose of the following *fundamental homomorphism theorem for modules*: let $f : M \rightarrow N$ be a homomorphism of left R -modules; then there is an induced R -module isomorphism $f' : M/\ker(f) \rightarrow \text{img}(f)$ given by $f'(x + \ker(f)) = f(x)$ for $x \in M$. Although this result can be deduced quickly from the fundamental homomorphism theorem for groups (see §1.7), we reprove the full result here for emphasis. The crucial first step is to check that f' is well-defined. Let $K = \ker(f)$, and suppose $x + K = z + K$ for some $x, z \in M$. Then $z = x + k$ for some $k \in K$, and

$$f'(z + K) = f(z) = f(x + k) = f(x) + f(k) = f(x) + 0 = f(x) = f'(x + K),$$

so that f' is well-defined. Second, let us check that f' is an R -module homomorphism. For all $x, y \in M$ and $r \in R$, compute

$$f'((x + K) + (y + K)) = f'((x + y) + K) = f(x + y) = f(x) + f(y) = f'(x + K) + f'(y + K);$$

$$f'(r(x + K)) = f'((rx) + K) = f(rx) = rf(x) = rf'(x + K).$$

Third, f' is injective, since for $x, y \in M$, $f'(x + K) = f'(y + K)$ implies $f(x) = f(y)$, hence $f(x - y) = 0$, hence $x - y \in K$, hence $x + K = y + K$. Fourth, the image of f' is $\{f'(x + K) : x \in M\} = \{f(x) : x \in M\}$, which is the image of f , so that f' is surjective as a mapping into the codomain $\text{img}(f)$.

The following generalization of the fundamental theorem is sometimes useful. Let $f : M \rightarrow N$ be an R -module homomorphism, and let H be any submodule of M . We would like to define an induced homomorphism $f' : M/H \rightarrow N$ by setting $f'(x + H) = f(x)$ for $x \in M$. When is this definition permissible? Answer: f' is well-defined iff $H \subseteq \ker(f)$. For in this case, $x + H = z + H$ (where $x, z \in M$) implies $z = x + h$ for some $h \in H$. Since $f(h) = 0$,

$$f'(z + H) = f(z) = f(x + h) = f(x) + f(h) = f(x) = f'(x + H).$$

Conversely, one may check that f' is multi-valued (i.e., not a well-defined function) if there exists $h \in H$ with $f(h) \neq 0$. In the case where f' is well-defined, it is routine to check (as above) that f' is an R -module homomorphism with $\text{img}(f') = \text{img}(f)$. Moreover, one may verify that $\ker(f')$ is precisely $\ker(f)/H = \{x + H : x \in \ker(f)\}$.

The generalized fundamental theorem can be restated as the following *universal mapping property (UMP) for the quotient module M/H* . Let $f : M \rightarrow N$ be an R -module homomorphism, let H be a submodule of M , and let $p : M \rightarrow M/H$ be the canonical projection. If $H \subseteq \ker(f)$, there exists a unique R -module homomorphism $f' : M/H \rightarrow N$

such that $f = f' \circ p$. Moreover, $\text{img}(f') = \text{img}(f)$ and $\ker(f') = \ker(f)/H$. The setup in this theorem can be visualized by the following diagram:

$$\begin{array}{ccc} M & \xrightarrow{p} & M/H \\ & \searrow f & \downarrow \exists! f' \\ & & N \end{array}$$

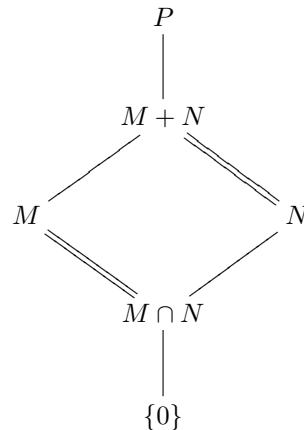
Existence of f' follows from the discussion above, together with the observation that $(f' \circ p)(x) = f'(p(x)) = f'(x + H) = f(x)$ for all $x \in M$. Similarly, uniqueness of f' holds since the requirement $f = f' \circ p$ forces us to define $f'(x + H) = f'(p(x)) = (f' \circ p)(x) = f(x)$ for all $x \in M$, as we did above.

17.9 More Module Homomorphism Theorems

We now use the fundamental homomorphism theorem for modules to deduce some further isomorphism theorems that play a prominent role in module theory.

Diamond Isomorphism Theorem. *Let M and N be submodules of a left R -module P . The quotient modules $M/(M \cap N)$ and $(M + N)/N$ are isomorphic, via the R -map $g : M/(M \cap N) \rightarrow (M + N)/N$ given by $g(m + M \cap N) = m + N$ for $m \in M$.* To prove this, apply the fundamental homomorphism theorem to the map $f : M \rightarrow (M + N)/N$ given by $f(m) = m + N$ for $m \in M$. The map f is the composition of an inclusion $M \rightarrow M + N$ and a canonical map $M + N \rightarrow (M + N)/N$, hence is a module homomorphism. Since, for all m in the domain M of f , $f(m) = 0$ in $(M + N)/N$ iff $m + N = 0 + N$ iff $m \in N$, the kernel of f is precisely $M \cap N$. We claim that the image of f is all of $(M + N)/N$. A typical element of this set is a coset $(m + n) + N$, where $m \in M$ and $n \in N$. But $(m + n) + N = (m + N) + (n + N) = (m + N) + 0 = m + N = f(m)$, so that this coset is in the image of f . Applying the fundamental theorem to f now provides an isomorphism g given by the stated formula.

The following picture, showing part of the submodule lattice of P , explains the name “diamond isomorphism theorem” and can aid in remembering the theorem:



The edges marked by double lines indicate which two quotient modules are isomorphic. Note that by interchanging M and N , we also have an R -module isomorphism $(M + N)/M \cong N/(M \cap N)$.

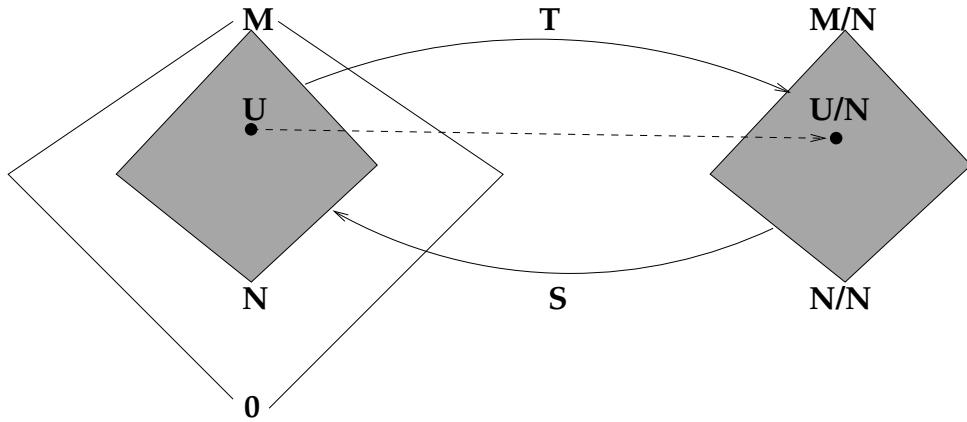
Nested Quotient Isomorphism Theorem. *Given a left R -module M with submodules H and N such that $H \subseteq N \subseteq M$, there is an R -module isomorphism $g : (M/H)/(N/H) \rightarrow (M/N)$ given by $g((x+H)+(N/H)) = x+N$ for $x \in M$.* To prove this, start with the projection homomorphism $p : M \rightarrow M/N$. Since $H \subseteq \ker(p) = N$ by assumption, applying the generalized fundamental homomorphism theorem to p produces a well-defined R -map $f : M/H \rightarrow M/N$ given by $f(x+H) = x+N$ for $x \in M$. The generalized theorem tells us that $\text{img}(f) = \text{img}(p) = M/N$ and $\ker(f) = \ker(p)/H = N/H$. Applying the fundamental homomorphism theorem to f , we get the isomorphism g appearing in the theorem statement.

Correspondence Theorem for Modules. *Let N be a fixed submodule of a left R -module M , and let $p : M \rightarrow M/N$ be the canonical map. Let A be the collection of all submodules U of M such that $N \subseteq U \subseteq M$. Let B be the collection of all submodules of M/N . There are inclusion-preserving, mutually inverse bijections $T : A \rightarrow B$ and $S : B \rightarrow A$ given by $T(U) = p[U] = U/N$ (the direct image of U under p) for $U \in A$ and $S(V) = p^{-1}[V]$ (the inverse image of V under p) for $V \in B$. Note that square brackets denote the direct or inverse images of subsets under p , whereas round parentheses denote ordinary function evaluation. Here, U is an element of the domain of T , whereas U is a subset of the domain of p ; similarly for V , S , and p^{-1} . In particular, $T \neq p$ as functions and $S \neq p^{-1}$; indeed, the function p^{-1} only exists when $N = \{0\}$. The statement that T preserves inclusions means that for all $U_1, U_2 \in A$, if $U_1 \subseteq U_2$ then $T(U_1) \subseteq T(U_2)$; similarly for S .*

The proof of the correspondence theorem consists of a sequence of routine but lengthy verifications, some details of which appear in Exercises 38 and 42. First, T does map A into B , since the direct image $p[U]$ is a submodule of M/N for any submodule U of M (whether or not U contains N). Second, S does map B into A , since the inverse image $p^{-1}[V]$ is a submodule of M containing $\ker(p) = N$ for any submodule V of M/N . Third, for any subset W of M/N , $p[p^{-1}[W]] = W$ follows from the surjectivity of the function $p : M \rightarrow M/N$. Consequently, taking W to be any submodule V of M/N , we see that $T(S(V)) = V = \text{id}_B(V)$ for all $V \in B$, and hence $T \circ S = \text{id}_B$. Fourth, let us show that $S(T(U)) = U = \text{id}_A(U)$ for any $U \in A$ (hence $S \circ T = \text{id}_A$). Set theory alone implies that $S(T(U)) = p^{-1}[p[U]] \supseteq U$, so it suffices to check the reverse inclusion. Recall that U is a submodule of M containing N . Let $x \in p^{-1}[p[U]]$. Then $p(x) = x + N \in p[U] = U/N$, so there exists $z \in U$ with $x + N = z + N$. In turn, there exists $n \in N$ with $x = z + n$. Since $N \subseteq U$ and U is closed under addition, we have $x \in U$. This establishes the inclusion $p^{-1}[p[U]] \subseteq U$. Fifth, one can show that inclusions are preserved when we take direct images or inverse images of subsets under any function. This fact (applied to p) implies that T and S preserve inclusions.

It follows that T and S are *lattice isomorphisms* between the lattice A of submodules of M containing N and the lattice B of all submodules of M/N . (A lattice isomorphism is a bijection f between two lattices such that f and its inverse preserve the underlying order relation, which in this instance is set inclusion.) So the correspondence theorem provides some retroactive motivation for the quotient module construction: if we are studying the submodule lattice of M , we can focus attention on the part of the lattice “above” the submodule N by passing to the submodule lattice of M/N . See Figure 17.2 for a picture of the relevant lattices.

Recognition Theorem for Direct Products. *Suppose N and P are submodules of a left R -module M such that $N+P = M$ and $N \cap P = \{0_M\}$. There is an R -module isomorphism $g : N \times P \rightarrow M$ given by $g((x,y)) = x+y$ for $x \in N$ and $y \in P$.* We prove this by checking that g is a one-to-one, onto, R -linear map. To verify R -linearity, fix $x_1, x_2 \in N$

**FIGURE 17.2**

The Lattice Isomorphisms in the Correspondence Theorem for Modules.

and $y_1, y_2 \in P$, and compute

$$\begin{aligned} g((x_1, y_1) + (x_2, y_2)) &= g((x_1 + x_2, y_1 + y_2)) = (x_1 + x_2) + (y_1 + y_2) \\ &= (x_1 + y_1) + (x_2 + y_2) = g((x_1, y_1)) + g((x_2, y_2)). \end{aligned}$$

Similarly, for $x \in N$, $y \in P$, and $r \in R$,

$$g(r(x, y)) = g((rx, ry)) = rx + ry = r(x + y) = rg((x, y)).$$

To see g is one-to-one, we show that $\ker(g) = \{(0, 0)\}$. For any (x, y) in $\ker(g)$, $0 = g((x, y)) = x + y$, so $y = -x$. We know $y \in P$, and since $y = -x$, y is also in the submodule N . Thus $y \in N \cap P = \{0_M\}$, hence $y = -x = 0$. Then $(x, y) = (0, 0)$ as needed. Finally, to see g is onto, let $z \in M$ be arbitrary. Since $N + P = M$, there exist $x \in N$ and $y \in P$ with $z = x + y = g((x, y))$.

In Exercise 49, we generalize the Recognition theorem to products of more than two factors.

17.10 Chains of Submodules

In Figure 17.1, all maximal paths through the submodule lattice from M to $\{0\}$ have the same length, namely three steps. Our next theorem shows this must always happen when there exists a maximal path through the submodule lattice of finite length.

Fix a ring R and a left R -module M . A *chain of submodules* of M is a list (M_0, M_1, \dots, M_m) where each M_i is a submodule of M and $M_0 \supsetneq M_1 \supsetneq \dots \supsetneq M_m$. We say that this chain has *length* m . Chains of submodules (M_0, M_1, \dots) of infinite length are defined similarly. A chain of finite length is called a *maximal chain* iff it cannot be extended to a longer chain by adding another submodule to the beginning, middle, or end of the list. Observe that the chain (M_0, \dots, M_m) is maximal iff $M_0 = M$ and $M_m = \{0\}$ and for all i between 1 and m , there exists no submodule P of M with $M_{i-1} \supsetneq P \supsetneq M_i$. By the correspondence theorem, the last condition is equivalent to saying that the only

submodules of the quotient module M_{i-1}/M_i are $\{0\} = M_i/M_i$ and M_{i-1}/M_i . In other words, the chain of submodules (M_0, \dots, M_m) is maximal iff $M = M_0$ and $M_m = \{0\}$ and M_{i-1}/M_i is a simple module for $1 \leq i \leq m$.

We can now state the **Jordan–Hölder Theorem for Modules**. *Let M be a module over a ring R such that there exists a maximal chain of submodules (M_0, M_1, \dots, M_m) of finite length m . Then: (a) any chain (N_0, N_1, \dots) of submodules of M has finite length $n \leq m$; (b) if the chain (N_0, \dots, N_n) is also maximal, then $n = m$ and for some permutation f of $\{1, 2, \dots, m\}$, we have $M_{i-1}/M_i \cong N_{f(i)-1}/N_{f(i)}$ for $1 \leq i \leq m$.*

We prove the theorem by induction on m , the length of the given maximal chain. If $m = 0$, then $M = M_0 = M_m = \{0\}$, and the conclusions of the theorem are evident. If $m = 1$, then $M \cong M/\{0\} = M_0/M_1$ must be a simple module, and again the needed conclusions follow immediately. Now assume $m > 1$ and that the theorem is already known for all R -modules having a maximal chain of length less than m . Fix a chain (N_0, N_1, \dots) of submodules of M . We first prove (a).

Case 1: $N_1 \subseteq M_1$. We can apply the induction hypothesis to the R -module M_1 , which has a maximal chain (M_1, \dots, M_m) of length $m - 1$. The chain (N_1, N_2, \dots) must therefore have some finite length $n - 1 \leq m - 1$, so that (N_0, N_1, \dots, N_n) has finite length $n \leq m$.

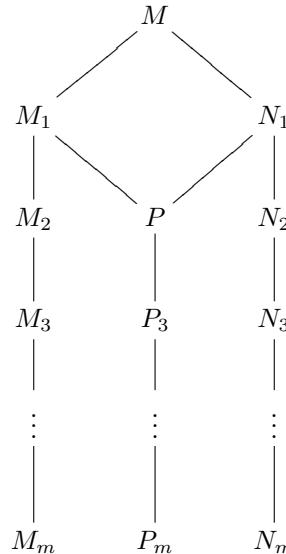
Case 2: N_1 is not a subset of M_1 . Note $M \neq N_1 \neq M_1$, but M and M_1 are the only two submodules of M containing M_1 . So M_1 is not a subset of N_1 , and $P = M_1 \cap N_1$ is a proper submodule of both M_1 and N_1 . Moreover, $M_1 + N_1$ is a submodule properly containing M_1 , so $M = M_1 + N_1$. By the diamond isomorphism theorem, $N_1/P = N_1/(N_1 \cap M_1) \cong (N_1 + M_1)/M_1 = M/M_1$ is a simple module. Starting with the chain (M_1, P) , we can build longer and longer chains of submodules of M_1 by repeatedly inserting a new submodule somewhere in the existing chain, as long as this can be done. By the induction hypothesis applied to M_1 , this insertion process must terminate in finitely many steps. When it does terminate, we must have by definition a maximal chain $(P_1 = M_1, P_2, \dots, P_m)$ of submodules of M_1 , which must have length $m - 1$ by induction applied to the module M_1 . We know $M_1 \cap N_1 = P = P_i$ for some i with $2 \leq i \leq m$. Since N_1/P is simple, $(N_1, P = P_i, P_{i+1}, \dots, P_m)$ is a maximal chain of submodules of N_1 of finite length $m - i + 1 \leq m - 1$. So the induction hypothesis is applicable to the R -module N_1 , and we conclude that the chain (N_1, N_2, \dots) of submodules of N_1 has finite length at most $m - i + 1 \leq m - 1$. So (N_0, N_1, \dots, N_n) has finite length $n \leq m$, proving (a) for the module M .

To prove (b), we now assume that (N_0, N_1, \dots, N_n) is a maximal chain of submodules of M . Using part (a) with the roles of the two maximal chains reversed, we immediately get $m \leq n$ and hence $n = m$. To obtain the further conclusion about the isomorphism of quotient modules, we again consider two cases.

Case 1: $N_1 \subseteq M_1$. Since $M_1 \subsetneq M$, maximality of the chain (N_0, N_1, \dots, N_n) forces $M_1 = N_1$. Then $M_0/M_1 = N_0/N_1$. We can match up the rest of the quotient modules by applying part (b) of the induction hypothesis to the module $M_1 = N_1$ and the maximal chains (M_1, \dots, M_m) and (N_1, \dots, N_n) within this module.

Case 2: N_1 is not a subset of M_1 . As before, we let $P = M_1 \cap N_1$ and note that $M = M_1 + N_1$. By the diamond isomorphism theorem and the assumed maximality of the two given chains, $N_1/P \cong M_0/M_1$ and $M_1/P \cong N_0/N_1$ are simple modules. As before, we can extend the chain (M_1, P) to a maximal chain $(M_1, P, P_3, \dots, P_m)$ of submodules of M_1 .

We now have four maximal chains of submodules of M that look like this:



Note that $(N_1, P, P_3, \dots, P_m)$ is also a maximal chain of submodules of N_1 . Applying the induction hypothesis to M_1 , we know that the modules $M_0/M_1, M_1/M_2, \dots, M_{m-1}/M_m$ are isomorphic (in some order) to the modules $M/M_1, M_1/P, P/P_3, \dots, P_{m-1}/P_m$. By the diamond isomorphisms at the top of the figure, these modules (switching the first two) are isomorphic to $M/N_1, N_1/P, P/P_3, \dots, P_{m-1}/P_m$. Applying the induction hypothesis to N_1 , we know that the modules just listed are isomorphic (in some order) to the modules $N_0/N_1, N_1/N_2, \dots, N_{m-1}/N_m$. Combining all these steps, the modules M_{i-1}/M_i (for $1 \leq i \leq m$) are isomorphic in some order to the modules N_{i-1}/N_i . This completes the proof of (b).

17.11 Modules of Finite Length

We say that an R -module M has *finite length* iff there exists a maximal chain of submodules of M of finite length. In this case, we have proved that all maximal chains of submodules of M have the same length m , and we define the *length of the module M* to be $\text{len}(M) = m$. If M does not have finite length, set $\text{len}(M) = \infty$.

Note that $\text{len}(M) = 0$ iff M is the zero module, whereas $\text{len}(M) = 1$ iff M is a simple module. One readily checks that if $M \cong M'$, then $\text{len}(M) = \text{len}(M')$. Next, suppose M is any R -module with a submodule N . We claim $\text{len}(M) < \infty$ iff $\text{len}(N) < \infty$ and $\text{len}(M/N) < \infty$, in which case $\text{len}(M) = \text{len}(N) + \text{len}(M/N)$. On one hand, suppose $\text{len}(M) = m < \infty$. Any chain of submodules of N is also a chain of submodules of M , so has length at most m . Hence N has finite length. On the other hand, by the correspondence theorem, any chain of submodules of M/N has the form $M_0/N \supsetneq M_1/N \supsetneq \dots \supsetneq M_k/N$ for some submodules $M_0 \supsetneq M_1 \supsetneq \dots \supsetneq M_k$ of M containing N . So the given chain in M/N has length at most m . Conversely, suppose $\text{len}(N) = n < \infty$ and $\text{len}(M/N) = k < \infty$. Fix a maximal chain of submodules of N , say $N = N_0 \supsetneq N_1 \supsetneq \dots \supsetneq N_n = \{0\}$. Similarly, fix a maximal chain of submodules of M/N , which must have the form $M/N = M_0/N \supsetneq M_1/N \supsetneq \dots \supsetneq M_k/N = N/N$ for certain

submodules $M = M_0 \supsetneq M_1 \supsetneq \cdots \supsetneq M_k = N$. Splicing these chains together, we get a chain of submodules

$$M = M_0 \supsetneq M_1 \supsetneq \cdots \supsetneq M_k = N = N_0 \supsetneq N_1 \supsetneq \cdots \supsetneq N_n = \{0\}$$

of length $k+n$. This chain must be maximal, since otherwise the chain for N or the chain for M/N would not have been maximal. This proves $\text{len}(M) = k+n = \text{len}(M/N)+\text{len}(N) < \infty$.

Given R -modules M_1, \dots, M_k , we show by induction that the product module $M = M_1 \times \cdots \times M_k$ has finite length iff all M_i have finite length, in which case $\text{len}(M) = \text{len}(M_1) + \cdots + \text{len}(M_k)$. This is evident if $k = 1$. Assume $k > 1$ and the result is known for products of $k-1$ modules. On one hand, if $\text{len}(M) < \infty$, then each submodule $\{0\} \times \cdots \times M_i \times \cdots \times \{0\}$ has finite length. As M_i is isomorphic to this submodule, $\text{len}(M_i) < \infty$. On the other hand, suppose each M_i has finite length. Note $M = N \times M_k$ where $N = M_1 \times \cdots \times M_{k-1}$. By induction, $\text{len}(N \times \{0\}) = \text{len}(N) = \text{len}(M_1) + \cdots + \text{len}(M_{k-1}) < \infty$. We have $M/(N \times \{0\}) \cong M_k$, which has finite length. By the result in the previous paragraph, $\text{len}(M) = \text{len}(N) + \text{len}(M_k) = \sum_{i=1}^k \text{len}(M_i) < \infty$.

17.12 Free Modules

To motivate our discussion of free modules, we first recall from §1.8 some fundamental definitions and theorems concerning bases of a vector space. Let V be a vector space over a field F . A subset S of V spans V iff every element of V can be written as a finite linear combination $a_1s_1 + \cdots + a_ks_k$ for some $k \in \mathbb{N}$, $a_i \in F$, and $s_i \in S$. An ordered list (s_1, \dots, s_k) of elements of V is *linearly independent* iff the relation $a_1s_1 + \cdots + a_ks_k = 0$ ($a_i \in F$) implies that every $a_i = 0$. A subset S of V (finite or not) is *linearly independent* iff every finite list of distinct elements of S is linearly independent. A subset S is a *basis* for V iff it is linearly independent and spans V . In this case, every $v \in V$ can be written *uniquely* as a linear combination of finitely many elements of S . In the finite-dimensional case, we usually use an (ordered) list rather than an (unordered) subset to specify a basis. We proved in Chapter 16 that every vector space V has a basis, and every basis of V has the same cardinality. If X is a basis of V , then every function from X into an F -vector space W has a unique extension to an F -linear map from V into W . Given any set X and field F , there exists an F -vector space having X as a basis. Every vector space is isomorphic to a direct sum of copies of F , as we saw in §17.5.

Now we consider the analogous concepts for a general left R -module M . The definitions are essentially the same, but some of the terminology changes. We have already discussed the notion of spanning: for $S \subseteq M$, M is generated by S iff $\langle S \rangle = M$ iff every element of M can be written as a finite sum $a_1s_1 + \cdots + a_ks_k$ for some $k \in \mathbb{N}$, $a_i \in R$, and $s_i \in S$. Recall that such a sum is called an *R -linear combination of elements of S* . Note $\langle \emptyset \rangle = \{0_M\}$. An ordered list (s_1, \dots, s_k) of elements of M is called *R -independent* (or *R -linearly independent*) iff the relation $a_1s_1 + \cdots + a_ks_k = 0$ ($a_i \in R$) implies that every $a_i = 0$. Negating this definition, we see that the list (s_1, \dots, s_k) is *R -dependent* or *R -linearly dependent* iff there exist $a_1, \dots, a_k \in R$ with $a_1s_1 + \cdots + a_ks_k = 0$ and some $a_j \neq 0$. A subset S of M (finite or not) is *R -independent* iff every finite list of distinct elements of S is linearly independent; otherwise, S is *R -dependent*. If there exists an R -independent generating set S for the module M , then M is called a *free R -module*, and S is called an *R -basis* or *basis* of M .

Certain results about generating sets, independence, and free R -modules are exactly like the corresponding vector space results, but other familiar theorems about vector spaces and bases are *false* for left R -modules. Here are the key facts in the module case:

1. *If M is a free left R -module with basis S , then every $y \in M$ can be uniquely expressed as a finite R -linear combination of elements of S .* The existence of such an expression follows from the fact that S generates M . For uniqueness, suppose that $y = a_1s_1 + \dots + a_ks_k = b_1s_1 + \dots + b_ks_k$ where $k \in \mathbb{N}$, $a_i, b_i \in R$, and the s_i are distinct elements of S . (We can assume that the same elements of S are used in both expressions, by adding new terms with zero coefficients if needed.) Subtracting gives $0 = \sum_{i=1}^k (a_i - b_i)s_i$, so that $a_i - b_i = 0$ and $a_i = b_i$ for all i by the R -independence of S .
2. **Universal mapping property (UMP) for free R -modules.** *Let M be a free left R -module with basis S . Given any left R -module N and any function $g : S \rightarrow N$, there exists a unique R -module homomorphism $g' : M \rightarrow N$ whose restriction to S equals g .*

$$\begin{array}{ccc} S & \xrightarrow{\subseteq} & M \\ & \searrow g & \downarrow \exists! g' \\ & & N \end{array}$$

To prove existence, take any $y \in M$, and write y uniquely as $y = \sum_{x \in S} a_x x$ where $a_x \in R$ and all but finitely many a_x are zero. Define $g'(y) = \sum_{x \in S} a_x g(x) \in N$. One checks that g' is R -linear and that $g'(x) = g(x)$ for $x \in S$. For uniqueness, suppose g'' is another R -linear map such that g'' extends g . Then $g' = g''$ follows because two R -linear maps that agree on the generating set S must be equal.

3. *Every free left R -module M is isomorphic to a direct sum of copies of R .* We imitate the proof of the corresponding vector space result (see §17.5). Let S be an R -basis for the free module M . Let $N = \bigoplus_{x \in S} R$. Each element of N is a function $g : S \rightarrow R$ that is nonzero for only finitely many elements of its domain S . Define a map $T : N \rightarrow M$ by setting $T(g) = \sum_{x \in S} g(x) \cdot x$, which makes sense because the sum has only a finite number of nonzero terms. Define a map $T' : M \rightarrow N$ as follows. Given $y \in M$, write y uniquely as $y = \sum_{x \in S} c_x x$ where only finitely many $c_x \in R$ are nonzero. Define $T'(y)$ to be the function given by $g(x) = c_x$ for $x \in S$. One can verify that T and T' are R -linear maps that are inverses of each other. We call $T'(y) = (c_x : x \in S)$ the *coordinates of y relative to the basis S* .
4. *Given any set S and any nonzero ring R , there exists a free left R -module with basis S .* We let $M = \bigoplus_{x \in S} R$, the direct sum of S copies of R . Elements of M are functions $g : S \rightarrow R$ with finite support. For each $s \in S$, we have an associated function $e_s : S \rightarrow R$ such that $e_s(x) = 0_R$ for all $x \neq s$ in S , and $e_s(s) = 1_R$. (Note that $0_R \neq 1_R$ in the nonzero ring R .) Each e_s belongs to M ; we show that $\{e_s : s \in S\}$ is a basis for M . Given any nonzero $g \in M$, let $\{s_1, \dots, s_k\}$ be the elements of S for which $g(s_i) \neq 0$. Then $g = \sum_{i=1}^k g(s_i)e_{s_i}$, as can be seen by evaluating both sides at each $s \in S$. Next, to show R -independence, suppose $\sum_{t \in S} a_t e_t = 0$ for some $a_t \in R$ (with all but finitely many a_t equal to zero). Evaluating this function at $s \in S$, we see that each $a_s = 0$. Finally, we change notation by replacing each element $e_s \in M$ by the corresponding element $s \in S$. Then M is a free R -module with basis S (before the notation change, we had a basis $\{e_s\}$ in bijective correspondence with S).

5. *Not every left R -module is free in general.* For example, consider $R = \mathbb{Z}$. A free R -module M is isomorphic to a direct sum of copies of \mathbb{Z} ; hence, M is either the zero module or M is infinite. On the other hand, any finite commutative group G is a \mathbb{Z} -module. Therefore, if G is not the zero group, then G is a non-free \mathbb{Z} -module.

More generally, suppose R is an infinite ring containing a left ideal I such that R/I is finite and nonzero. Then R/I is a left R -module that cannot be isomorphic to a direct sum of copies of R ; hence R/I is a non-free R -module.

For certain rings R , every R -module *is* free. For example, this holds if R is a field (every vector space has a basis), or if R is the zero ring (here, $0_R = 1_R$ forces every R -module to be the zero module). The arguments in §16.9 prove that for a *division ring* R (i.e., a possibly non-commutative ring in which every nonzero element has a multiplicative inverse), every R -module is free.

6. *Not every left R -module is isomorphic to a submodule of a free R -module.* For example, consider $R = \mathbb{Z}$. If G is a nonzero finite commutative group, then G is not isomorphic to a submodule of any free R -module M . For, M must be a direct sum of copies of \mathbb{Z} , and therefore the only element of M of finite order is the identity element.

7. *Every left R -module is isomorphic to a quotient module of a free left R -module.*

Let M be an arbitrary left R -module. Let F be a free left R -module having the set M as a basis. Consider the identity map $\text{id}_M : M \rightarrow M$, where the domain is viewed as a subset of F , and the codomain is viewed as a left R -module. By the UMP, this map extends uniquely to an R -module homomorphism $g : F \rightarrow M$, which is evidently surjective. By the fundamental homomorphism theorem, $F/\ker(g) \cong M$. More generally, if S is any generating set for M , it suffices to let F be a free left R -module with S as a basis. The R -map $g : F \rightarrow M$ that extends the inclusion map $i : S \rightarrow M$ must be surjective, since the image of g is a submodule of M containing all the generators S of M . This shows that *every finitely generated left R -module M is isomorphic to a quotient of a finitely generated free left R -module F .*

17.13 Size of a Basis of a Free Module

One of the most frequently used facts about vector spaces over a field is that every vector space has a *unique* dimension. We have already seen that, for general rings, not every left R -module has a “dimension” (since not every left R -module is free). Even for R -modules M that are free, there may exist two bases of M of different cardinalities. We give an example of this phenomenon below.

However, the situation is more pleasant for commutative rings. We now prove that *if R is a nonzero commutative ring and N is a free left R -module with two bases X and Y , then $|X| = |Y|$.* We call the size of any basis of N the *dimension* of N . For the proof, we first consider the case where X and Y are finite ordered bases of N , say $X = (x_1, \dots, x_n)$ and $Y = (y_1, \dots, y_m)$. As in §17.7, we assume known the theorems from ring theory that R has a maximal ideal I and that $F = R/I$ is a field. We have seen that $V = N/IN$ is an R/I -module, i.e., an F -vector space. Let $X' = (x_1 + IN, \dots, x_n + IN)$ be the list of images of elements of X in the quotient module. We claim that X' is an ordered basis for the F -vector space V . One may check that X' spans (generates) V , since X generates N . To

show that X' is linearly independent, suppose $\sum_{i=1}^n (s_i + I)(x_i + IN) = 0$ for some $s_i \in R$. We must show that each $s_i + I = 0$ in R/I , i.e., that $s_i \in I$ for all i . From the given relation and the definition of the operations in N/IN , we deduce $(\sum_{i=1}^n s_i x_i) + IN = 0$. Therefore, the element $x = \sum_{i=1}^n s_i x_i$ lies in IN . By definition of IN , we can write $x = \sum_{j=1}^p a_j z_j$ where $a_j \in I$ and $z_j \in N$. Writing $z_j = \sum_{i=1}^n c_{i,j} x_i$ (for some $c_{i,j} \in R$) and substituting, we see that $x = \sum_{i=1}^n (\sum_{j=1}^p a_j c_{i,j}) x_i$. By independence of the x_i 's in N , we conclude that $s_i = \sum_{j=1}^p a_j c_{i,j}$ for all i . Since each a_j lies in the ideal I , we have $s_i \in I$ as needed. (Note that this step might fail if I were merely a left ideal in a non-commutative ring R .) So X' is a basis for V . By the same argument, $Y' = (y_1 + IN, \dots, y_m + IN)$ is an ordered basis for V . We now quote the theorem for vector spaces that says that the number of elements in a basis for V is unique (see §16.9). Since X' has length n and Y' has length m , that theorem gives $n = m$, as needed.

Now consider the general case, where X and Y are sets (possibly infinite). Let $p : N \rightarrow N/IN$ be the canonical map, and set $X' = p[X]$, $Y' = p[Y]$. As before, X' and Y' generate the F -vector space N/IN . Moreover, the preceding argument, applied to all finite lists of distinct elements of X , shows that X' (and similarly Y') are F -linearly independent subsets of N/IN . Therefore, $|X'| = |Y'|$ by the theorem for vector spaces. It now suffices to show that the restriction of p to X is injective, so that $|X| = |X'|$ (and similarly $|Y| = |Y'|$). This fact already follows from the argument used to show linear independence of X' . For if $x_1 \neq x_2$ in X , then the list (x_1, x_2) is R -independent in N , so that $(x_1 + IN, x_2 + IN)$ is F -independent in N/IN , so in particular $p(x_1) = x_1 + IN \neq x_2 + IN = p(x_2)$ in N/IN .

Next we give the promised example of a non-commutative ring R and a free R -module that has multiple R -bases with different cardinalities. Let F be any field, let $X = \{x_n : n \geq 0\}$ be a countably infinite set, and let V be an F -vector space with basis X . Let $R = \text{Hom}_F(V, V)$ be the set of linear transformations from V to itself. We have seen that R is an F -vector space. In fact, R is also a non-commutative ring, if we define multiplication of elements $f, g \in R$ to be composition of functions. The ring axioms may be routinely verified; in particular, the distributive laws follow since addition of functions is defined pointwise.

Like any ring, R is a left R -module. For each integer $k \geq 1$, we will produce an R -basis for R consisting of k elements. Fix $k \in \mathbb{N}^+$. We define certain elements f_j and g_j in R (for $0 \leq j < k$) by specifying how these linear maps operate on the basis X of V . Recall (by integer division) that every integer $n \geq 0$ can be written uniquely in the form $n = ki + j$, for some integers i, j with $0 \leq j < k$. Let f_j send x_{ki+j} to x_i for all $i \geq 0$; let f_j send all other elements of X to 0. Let g_j send x_i to x_{ki+j} for all $i \geq 0$. We have $f_j g_j = \text{id}_V$ for all j , since both sides have the same effect on the basis X . For the same reason, $f_{j'} g_j = 0_R$ for $j \neq j'$ between 0 and $k - 1$, and

$$g_0 f_0 + g_1 f_1 + \cdots + g_{k-1} f_{k-1} = \text{id}_V. \quad (17.1)$$

The last identity follows because, given any $x_n \in X$, we can write $n = ki + j$ for a unique j between 0 and $k - 1$. Then $g_j f_j(x_n) = x_n$ for this j , while $g_{j'} f_{j'}(x_n) = 0$ for all other $j' \neq j$.

We can now show that $B_k = (f_0, \dots, f_{k-1})$ is a k -element ordered R -basis for the left R -module R . Suppose f is any element of R . Multiplying (17.1) on the left by f , we get

$$(fg_0)f_0 + (fg_1)f_1 + \cdots + (fg_{k-1})f_{k-1} = f$$

where $fg_j \in R$. This identity shows that B_k is a generating set for the left R -module R . Next, to test R -independence, suppose

$$h_0 f_0 + \cdots + h_{k-1} f_{k-1} = 0_R$$

for some $h_j \in R$. For each j_0 between 0 and $k - 1$, multiply this equation on the right by g_{j_0} to obtain $h_{j_0} = 0_R$ (using the relations above to simplify products $f_j g_{j_0}$). Thus, B_k is an R -independent list. Note how prominently the non-commutativity of R entered into this proof.

We have now proved that B_k is an ordered basis for the left R -module R for all $k \geq 1$. So R has R -bases of every finite cardinality. On the other hand, viewing R as an F -module (i.e., as a vector space over F), we know that every F -basis of R must have the same (infinite) cardinality.

17.14 Summary

Let R be a ring. Here we summarize the definitions and results for R -modules, module homomorphisms, and free modules that were covered in this chapter.

Definitions

1. A *left R -module* is an additive commutative group M and a scalar multiplication $\cdot : R \times M \rightarrow M$ satisfying closure, left associativity, the two distributive laws, and the identity axiom.
2. *Right R -modules* are defined like left modules, using a multiplication $\cdot : M \times R \rightarrow M$ obeying right associativity; the two types of modules are equivalent for commutative rings.
3. A *homomorphism* of R -modules (left or right) is a map between modules that preserves addition and scalar multiplication. Homomorphisms of R -modules are also called *R -maps* or *R -linear maps*.
4. For fields F , F -modules and F -maps are the same as F -vector spaces and linear transformations. For $R = \mathbb{Z}$, \mathbb{Z} -modules and \mathbb{Z} -maps are the same as commutative groups and group homomorphisms.
5. A *submodule* of a module is a subset containing 0 and closed under addition and scalar multiplication (and hence under additive inverses).
6. Let $f : M \rightarrow N$ be an R -module homomorphism. The set of all $x \in M$ such that $f(x) = 0_N$ is the *kernel* of f ; the set of all $y \in N$ of the form $y = f(x)$ for some $x \in M$ is the *image* of f . These are submodules of M and N , respectively.
7. If S is a subset of a module M , an *R -linear combination* of elements of S is a finite sum $\sum_i a_i s_i$ where $a_i \in R$ and $s_i \in S$. The set of all such R -linear combinations is $\langle S \rangle$, the submodule *generated by S* . If $\langle S \rangle = M$, we say that S *spans M* or *generates M* , and that M is generated by S . If this holds for some finite set S , we say M is *finitely generated*.
8. A left R -module M is *cyclic* iff there exists $x \in M$ with $M = \langle x \rangle = Rx = \{rx : r \in R\}$. M is *simple* iff $M \neq \{0\}$ and the only submodules of M are $\{0\}$ and M .
9. A subset S of a module M is called *R -independent* iff for all $k \in \mathbb{N}$, $a_i \in R$, and distinct $s_i \in S$, $\sum_{i=1}^k a_i s_i = 0$ implies every $a_i = 0$.
10. An R -independent generating set of an R -module M is called an *R -basis* for M .

If M has an R -basis, M is called a *free R -module*. For R commutative, the size of any R -basis of a free R -module M is the *dimension* of M .

11. The *length* $\text{len}(M)$ of an R -module M is the maximum n such that there is a chain of submodules $M_0 \supseteq M_1 \supseteq \cdots \supseteq M_n$ of M , or ∞ if there is no such n .

Module constructions

1. *Submodules*: Every submodule of an R -module is itself an R -module. Intersections and sums of submodules are submodules. The direct or inverse image of a submodule under an R -map is a submodule.
2. *Submodule generated by a set*: If S is any subset of an R -module M , there exists a smallest submodule $\langle S \rangle$ of M containing S . This submodule can be defined as the intersection of all submodules of M containing S , or as the set of all R -linear combinations of elements of S (including zero).
3. *Direct products and direct sums*: If M_i is a left R -module for each i in a set I , then the set of functions f with domain I such that $f(i) \in M_i$ for all $i \in I$ is a left R -module under pointwise operations on functions. The subset of functions that are zero for all but finitely many $i \in I$ is a submodule. These modules are called the *direct product* and *direct sum* of the M_i 's and are denoted by $\prod_{i \in I} M_i$ and $\bigoplus_{i \in I} M_i$. A common special case is R^n , the module of n -tuples of elements of R , which is a free R -module having an n -element R -basis.
4. *Hom modules*: Let M and N be left R -modules. The set $\text{Hom}(M, N)$ of group homomorphisms from M to N is a left R -module. If R is commutative, the set $\text{Hom}_R(M, N)$ of R -module homomorphisms from M to N is a left R -module.
5. *Quotient Modules*: Let N be a submodule of an R -module M . The *quotient module* M/N consists of all cosets $x + N$ with $x \in M$. For $x, z \in M$, $x + N = z + N$ iff $x - z \in N$. The operations in the quotient module are defined by $(x + N) + (y + N) = (x + y) + N$ and $r \cdot (x + N) = (rx) + N$ for $x, y \in M$ and $r \in R$. The canonical projection $p : M \rightarrow M/N$, given by $p(x) = x + N$ for $x \in M$, is a surjective R -map with kernel N . If S generates M , then $p[S]$ generates M/N .
6. *Change of Scalars*: If T is a subring of R , any R -module can be regarded as a T -module. If $f : S \rightarrow R$ is a ring homomorphism, then an R -module M becomes an S -module via $s \star x = f(s) \cdot x$ for $s \in S$ and $x \in M$. If I is an ideal of R annihilating an R -module M (i.e., $ix = 0$ for all $i \in I$ and $x \in M$), then M can be regarded as an R/I -module via $(r + I) \bullet x = r \cdot x$ for $r \in R$ and $x \in M$. In particular, for any R -module N and any ideal I of R , N/IN is an R/I -module. Taking I to be a maximal ideal in a commutative ring R , we can convert R -modules to R/I -vector spaces.

Results about module homomorphisms and submodules

Let M and N be left R -modules.

1. If $f, g : M \rightarrow N$ are two R -maps agreeing on a generating set for M , then $f = g$.
2. *Fundamental Homomorphism Theorem for R -modules*: Let $f : M \rightarrow N$ be an R -module homomorphism. Then f induces an R -module isomorphism $f' : M/\ker(f) \rightarrow \text{img}(f)$ given by $f'(x + \ker(f)) = f(x)$ for $x \in M$.

3. *Universal Mapping Property for the Quotient Module M/H :* Let $f : M \rightarrow N$ be an R -module homomorphism, let H be a submodule of M , and let $p : M \rightarrow M/H$ be the canonical map. If $H \subseteq \ker(f)$, then there exists a unique R -module homomorphism $f' : M/H \rightarrow N$ such that $f = f' \circ p$. Moreover, $\text{img}(f') = \text{img}(f)$ and $\ker(f') = \ker(f)/H$.
4. *Diamond Isomorphism Theorem:* Let M and N be submodules of an R -module P . The R -map $g : M/(M \cap N) \rightarrow (M+N)/N$ given by $g(m+M \cap N) = m+N$ (for $m \in M$) is an R -module isomorphism.
5. *Nested Quotient Isomorphism Theorem:* Given $H \subseteq N \subseteq M$ with H and N submodules of M , there is an R -module isomorphism $g : (M/H)/(N/H) \rightarrow (M/N)$ given by $g((x+H)+(N/H)) = x+N$ for $x \in M$.
6. *Correspondence Theorem for Modules:* Assume N is a submodule of M , and let $p : M \rightarrow M/N$ be the canonical map. Let A be the collection of all submodules U of M such that $N \subseteq U \subseteq M$. Let B be the collection of all submodules of M/N . There are inclusion-preserving, mutually inverse bijections $T : A \rightarrow B$ and $S : B \rightarrow A$ given by $T(U) = p[U] = U/N$ (the direct image of U under p) and $S(V) = p^{-1}[V]$ (the inverse image of V under p).
7. *Recognition Theorem for Direct Products:* Let N and P be submodules of M such that $N + P = M$ and $N \cap P = \{0_M\}$. There is an R -module isomorphism $g : N \times P \rightarrow M$ given by $g((x, y)) = x + y$ for $x \in N$ and $y \in P$.
8. *Jordan–Hölder Theorem for Modules:* If M has one maximal chain of submodules of finite length m , then every chain of submodules of M has length at most m , and every maximal chain has length m . Given two maximal chains $M_0 \supsetneq M_1 \supsetneq \dots \supsetneq M_m$ and $N_0 \supsetneq N_1 \supsetneq \dots \supsetneq N_m$ of M , the quotient modules $M_0/M_1, M_1/M_2, \dots, M_{m-1}/M_m$ are simple and are isomorphic (in some order) to $N_0/N_1, N_1/N_2, \dots, N_{m-1}/N_m$.
9. *Results for Finite Length Modules:* $M = \{0\}$ iff $\text{len}(M) = 0$. M is simple iff $\text{len}(M) = 1$. Isomorphic modules have the same length. For a submodule N of M , M has finite length iff N and M/N have finite length, and then $\text{len}(M) = \text{len}(N) + \text{len}(M/N)$. For $M = M_1 \times \dots \times M_k$, M has finite length iff all M_i have finite length, and then $\text{len}(M) = \text{len}(M_1) + \dots + \text{len}(M_k)$.

Results about free R -modules

1. If M is a free R -module with basis S , then every $y \in M$ can be uniquely expressed as a (finite) R -linear combination of elements of S .
2. *Universal mapping property for free R -modules:* Let M be a free R -module with basis X . Given any R -module N and any function $g : X \rightarrow N$, there exists a unique R -module homomorphism $g' : M \rightarrow N$ whose restriction to X equals g .
3. Every free R -module M is isomorphic to a direct sum of copies of R . Conversely, any such direct sum is a free R -module.
4. Given any set X and any nonzero ring R , there exists a free R -module with R -basis X .
5. Not every R -module is a free R -module. But if R is a field or division ring, every R -module is free.
6. Not every R -module is isomorphic to a submodule of a free R -module. But every R -module M is isomorphic to a quotient module of a free R -module F . If M is

finitely generated, F can be chosen to be finitely generated with the same number of generators as M .

7. If R is a nonzero commutative ring and N is a free R -module, then any two bases for N have the same cardinality. However, there exists a non-commutative ring R and a free R -module that has multiple R -bases with different cardinalities.
-

17.15 Exercises

Unless otherwise specified, assume R is an arbitrary ring in these exercises.

1. Let R be any nonzero ring. For each commutative group $(M, +)$ below, show that M is *not* a left R -module under the indicated scalar multiplication $\cdot : R \times M \rightarrow M$ by pointing out one or more module axioms that fail to hold.
 - (a) $M = R$, $r \cdot m = 0_M$ for all $r \in R$ and $m \in M$.
 - (b) $M = R$, $r \cdot m = m$ for all $r \in R$ and $m \in M$.
 - (c) $M = R$, $r \cdot m = r$ for all $r \in R$ and $m \in M$.
 - (d) $M = R^2$, $r \cdot (m_1, m_2) = (rm_1, m_2)$ for all $r, m_1, m_2 \in R$.
 - (e) $M = R$, $r \cdot m = mr$ for all $r, m \in R$ (assume R is non-commutative).
2. (a) For each $n \in \mathbb{N}^+$, show that the additive group R^n (viewed as a set of column vectors) is a left $M_n(R)$ -module if scalar multiplication $A \cdot v$ (for $A \in M_n(R)$ and $v \in R^n$) is defined to be the matrix-vector product Av .
 - (b) In (a), if we define $v \star A = Av$, do we get a right $M_n(R)$ -module structure on R^n ? Prove or give a counterexample.
 - (c) Show that R^n (viewed as a set of row vectors) is a right $M_n(R)$ -module if scalar multiplication $w \cdot A$ (for $A \in M_n(R)$ and $w \in R^n$) is defined to be the vector-matrix product wa .
3. Let V be a vector space over a field F , and let R be the ring of all F -linear transformations $T : V \rightarrow V$. Show that $(V, +)$ is a left R -module if we define $T \cdot v = T(v)$ for $T \in R$ and $v \in V$.
4. Let $(M, +)$ be a commutative group, and let $\cdot : R \times M \rightarrow M$ be a function. Define $\star : M \times R \rightarrow M$ by $m \star r = r \cdot m$ for all $r \in R$ and $m \in M$.
 - (a) Prove that each axiom (M1), (M2), (D1), and (D2) in §17.1 holds for \cdot iff the analogous axiom (M1'), (M2'), (D1'), (D2') holds for \star .
 - (b) Prove that, if R is commutative, then axiom (M3) holds for \cdot iff axiom (M3') holds for \star .
5. Let M, N be left R -modules and assume $f : M \rightarrow N$ is an R -linear map. Prove the following facts, which were stated in §17.1.
 - (a) The identity map $\text{id}_M : M \rightarrow M$ is an automorphism.
 - (b) The composition of R -maps is an R -map (similarly for monomorphisms, epimorphisms, isomorphisms, endomorphisms, and automorphisms).
 - (c) If f is an isomorphism, so is f^{-1} .
 - (d) Given a set X of left R -modules, the relation defined on X by $M \cong N$ iff M and N are isomorphic R -modules (for $M, N \in X$) is an equivalence relation.
 - (e) $f(0_M) = 0_N$ and $f(-x) = -f(x)$ for all $x \in M$.
 - (f) $0_R \cdot x = 0_M$ for all $x \in M$.
 - (g) $r \cdot 0_M = 0_M$ for all $r \in R$.
 - (h) $r \cdot (-x) = (-r) \cdot x = -(r \cdot x)$ for all $r \in R$ and all $x \in M$.
6. The complex number system \mathbb{C} is a ring, hence also a left \mathbb{C} -module and a left \mathbb{R} -module.
 - (a) Show that $f : \mathbb{C} \rightarrow \mathbb{C}$ defined by $f(a + ib) = a - ib$ for $a, b \in \mathbb{R}$ is a ring homomorphism and an \mathbb{R} -linear map, but not a \mathbb{C} -linear map.
 - (b) Show that $g : \mathbb{C} \rightarrow \mathbb{C}$ defined by $g(z) = iz$ for $z \in \mathbb{C}$ is a \mathbb{C} -linear map, but not a ring homomorphism.
 - (c) True or false: there exists **exactly one** function $\cdot : \mathbb{C} \times \mathbb{C} \rightarrow \mathbb{C}$ such that $(\mathbb{C}, +, \cdot)$ is a left \mathbb{C} -module.

7. *Opposite Rings.* Let $(R, +, \bullet)$ be a ring. Define an operation $\star : R \times R \rightarrow R$ by setting $a \star b = b \bullet a$ for all $a, b \in R$. (a) Prove that $(R, +, \star)$ is a ring. This ring is called the *opposite ring* to R and is denoted R^{op} . (b) Let $(M, +)$ be a commutative group. Show that $\cdot : M \times R \rightarrow M$ satisfies the axioms for a right R -module iff $* : R^{op} \times M \rightarrow M$, defined by $r * m = m \cdot r$ for $r \in R$ and $m \in M$, satisfies the axioms for a left R^{op} -module. (This result lets us reduce the study of right modules to the study of left modules over the opposite ring.)
8. Define a left R -module structure on the set of matrices $M_{m,n}(R)$. Is this a free R -module? If so, describe an R -basis.
9. Let $(M, +)$ be a commutative group. (a) Prove there is at most one scalar multiplication $\cdot : \mathbb{Z} \times M \rightarrow M$ that makes M a left \mathbb{Z} -module. (b) Prove that the scalar multiplication $\cdot : \mathbb{Z} \times M \rightarrow M$, defined in §17.2, does satisfy the axioms for a left \mathbb{Z} -module.
10. Let $(M, +)$ be a commutative group, and let $n \in \mathbb{N}^+$. (a) Prove there is at most one scalar multiplication $\cdot : \mathbb{Z}_n \times M \rightarrow M$ that makes M a left \mathbb{Z}_n -module. (b) Find and prove a condition on M that is necessary and sufficient for there to exist a left \mathbb{Z}_n -module with underlying additive group $(M, +)$.
11. Given a commutative group $(M, +)$, recall from Exercise 48 of Chapter 1 that we have the *endomorphism ring* $\text{End}(M)$ consisting of all group homomorphisms $f : M \rightarrow M$. We add and multiply $f, g \in \text{End}(M)$ via the rules $(f + g)(x) = f(x) + g(x)$ and $(f \circ g)(x) = f(g(x))$ for all $x \in M$. Prove: for any subring S of $\text{End}(M)$, M is a left S -module under the scalar multiplication $f \cdot x = f(x)$ for $f \in S$ and $x \in M$.
12. Let $(M, +)$ be a fixed commutative group with endomorphism ring $\text{End}(M)$.
 - (a) Given a scalar multiplication $\cdot : R \times M \rightarrow M$ satisfying the axioms for a left R -module, define a map $L_r : M \rightarrow M$ (for each $r \in R$) by setting $L_r(x) = r \cdot x$ for $x \in M$. The map L_r is called *left multiplication by r* . Confirm that $L_r \in \text{End}(M)$ for each $r \in R$. Then show that the map $L : R \rightarrow \text{End}(M)$, defined by $L(r) = L_r$ for $r \in R$, is a ring homomorphism. (b) Conversely, suppose $T : R \rightarrow \text{End}(M)$ is a given ring homomorphism. Define $\star : R \times M \rightarrow M$ by $r \star x = T(r)(x)$ for all $r \in R$ and all $x \in M$. Show that this scalar multiplication map turns M into a left R -module. (c) Let X be the set of all scalar multiplication functions $\cdot : R \times M \rightarrow M$ satisfying the axioms for a left R -module, and let Y be the set of all ring homomorphisms $L : R \rightarrow \text{End}(M)$. The constructions in (a) and (b) define maps $\phi : X \rightarrow Y$ and $\psi : Y \rightarrow X$. Check that $\phi \circ \psi = \text{id}_Y$ and $\psi \circ \phi = \text{id}_X$, so that ϕ and ψ are bijections. (This problem shows that “left R -module structures” on a given commutative group $(M, +)$ correspond bijectively with ring homomorphisms of R into the endomorphism ring $\text{End}(M)$.)
13. (a) Prove that for any ring S , there exists exactly one ring homomorphism $f : \mathbb{Z} \rightarrow S$. (b) Use (a) and Exercise 12 to give a very short proof that for every commutative group $(M, +)$, there exists a unique scalar multiplication that turns M into a left \mathbb{Z} -module.
14. Show that N is a submodule of an R -module M iff N is a nonempty subset of M closed under subtraction and left multiplication by scalars in R .
15. (a) Carefully check that every additive subgroup of a \mathbb{Z} -module M is automatically a \mathbb{Z} -submodule. (b) Give a specific example of a ring R , a left R -module M , and an additive subgroup N of M that is not a submodule of M . (c) Prove or disprove: for all $n \in \mathbb{N}^+$, every additive subgroup of a \mathbb{Z}_n -module M is automatically a \mathbb{Z}_n -submodule.

16. (a) Prove or disprove: if M and N are submodules of a left R -module P , then $M \cup N$ is a submodule of P . (b) Prove: if $\{M_i : i \in I\}$ is an indexed family of submodules of a left R -module P such that $I \neq \emptyset$ and for all $i, j \in I$, either $M_i \subseteq M_j$ or $M_j \subseteq M_i$, then $N = \bigcup_{i \in I} M_i$ is a submodule of P .
17. (a) Given submodules M and N of a left R -module P , confirm that $M + N$ is a submodule of P . (b) Given a family $\{M_i : i \in I\}$ of submodules of an R -module P , confirm that $\sum_{i \in I} M_i$ is a submodule of P . (c) In (b), prove that if N is any R -submodule of P containing every M_i , then $\sum_{i \in I} M_i \subseteq N$ (so the sum of the M_i is the smallest submodule containing every M_i).
18. Let S be any set, and let X be the set of all subsets of S . X becomes a partially ordered set via $U \leq V$ iff $U \subseteq V$ (for $U, V \in X$). Prove that X is a complete lattice.
19. (a) Prove that every simple R -module is cyclic. Give an example to show that the converse is not true in general. (b) Show that a \mathbb{Z} -module M is simple iff $|M|$ is prime. (You may need Lagrange's theorem from group theory.) (c) For a field F , show that an F -module M is simple iff $\dim_F(M) = 1$. (d) For a commutative ring R , show that the left R -module R is simple iff R is a field.
20. Give an example (with proof) of an infinite ring R and an infinite cyclic R -module M such that there exists a *unique* $x \in M$ with $M = Rx$.
21. (a) For any left ideal I of R , prove R/I is a cyclic left R -module. (b) Conversely, prove that every cyclic left R -module M is isomorphic to a module R/I for some left ideal I of R . (Use the fundamental homomorphism theorem.)
22. A subset I of R is called a *maximal left ideal* iff I is a submodule of the left R -module R such that $I \neq R$ and for any submodule J with $I \subseteq J \subseteq R$, either $J = I$ or $J = R$. (a) Prove that if I is a maximal left ideal of R , then R/I is a simple left R -module. (Use the correspondence theorem.) (b) Conversely, prove that every simple left R -module M is isomorphic to a module R/I for some maximal left ideal I of R .
23. We know that the set $M_2(\mathbb{R})$ of 2×2 real matrices is an additive group, a ring, a left \mathbb{Z} -module, a left \mathbb{R} -module, and a left $M_2(\mathbb{R})$ -module. Let N be the set of matrices of the form $\begin{bmatrix} 0 & a \\ 0 & b \end{bmatrix}$ for some $a, b \in \mathbb{R}$. (a) Show that N is an \mathbb{R} -submodule of $M_2(\mathbb{R})$, and an $M_2(\mathbb{R})$ -submodule of $M_2(\mathbb{R})$, but not an ideal of $M_2(\mathbb{R})$. (b) Show that N is *not* a simple \mathbb{R} -module. (c) Show that N *is* a simple $M_2(\mathbb{R})$ -module. (Show that any nonzero $M_2(\mathbb{R})$ -submodule P of N must be equal to N .)
24. (a) Give an example of submodules A, B, C of the \mathbb{R} -module $\mathbb{R} \times \mathbb{R}$ such that the “distributive law” $A \cap (B + C) = (A \cap B) + (A \cap C)$ is **not** true. (b) Let A, B , and C be submodules of a left R -module M . Prove: if $A \subseteq C$, then $A + (B \cap C) = (A + B) \cap C$.
25. Let S be a subset of a left R -module M . (a) Show that the set N' of R -linear combinations of elements of S is an R -submodule of M by checking the closure conditions in the definition. (b) Show that N' is a submodule of M by verifying that $N' = \sum_{s \in S} Rs$.
26. Let N be a subset of a left R -module M . Let I be the set of all $r \in R$ such that $r \cdot x = 0_M$ for all $x \in N$. (a) Prove that I is a submodule of the left R -module R . (b) Now assume N is a submodule of M . Prove that I is a (two-sided) ideal in the ring R .

27. Assuming $R \neq \{0\}$, prove that $R[x]$ is not a finitely generated R -module.
28. (a) Given an index set I and left R -modules M_i for each $i \in I$, verify that $N = \prod_{i \in I} M_i$ satisfies all the module axioms. (b) Check that N_0 , the set of $f \in N$ with finite support, is a submodule of N .
29. Verify that the maps S and T in §17.5 are R -linear maps that are inverses of each other.
30. Let M and N be left R -modules. Prove that $\text{Hom}_R(M, N)$ is always an additive subgroup of N^M .
31. *Bimodules.* Given rings R and S , an R, S -bimodule is a commutative group $(M, +)$ that has both a left R -module structure, given by $\cdot : R \times M \rightarrow M$, and a right S -module structure, given by $\star : M \times S \rightarrow M$, that are connected by the axiom $(r \cdot x) \star s = r \cdot (x \star s)$ for all $r \in R$, $s \in S$, and $x \in M$. (a) Prove that any ring R is an R, R -bimodule if we take \cdot and \star to be the multiplication operation of R . (b) Prove that R^n (viewed as column vectors) is an $M_n(R), R$ -bimodule using the natural action of matrices and scalars on column vectors. (c) Let M be a left R -module and N an R, S -bimodule. Show that the commutative group $\text{Hom}_R(M, N)$ of R -maps from M to N is a right S -module if we define $f \cdot s$ (for $f \in \text{Hom}_R(M, N)$ and $s \in S$) to be the function from M to N sending $x \in M$ to $f(x) \star s$. (d) Let M be an R, S -bimodule and N a left R -module. Show that the commutative group $\text{Hom}_R(M, N)$ of R -maps from M to N is a left S -module if we define $s \cdot f$ (for $f \in \text{Hom}_R(M, N)$ and $s \in S$) to be the function from M to N sending $x \in M$ to $f(x \star s)$.
32. Let M be a left R -module with submodule N . Define a relation \equiv on M by setting $x \equiv y$ iff $x - y \in N$ (for all $x, y \in M$). Prove \equiv is an equivalence relation, and prove the equivalence class of x is the coset $x + N$.
33. (a) Prove that $R[x]$ is a free left R -module by finding an explicit basis for this module. (b) Assume R is an integral domain and $g \in R[x]$ is monic of degree $n > 0$. Let $I = R[x]g$. Prove $R[x]/I$ is a free R -module with ordered basis $(1 + I, x + I, x^2 + I, \dots, x^{n-1} + I)$.
34. Suppose S is a ring, $f : S \rightarrow R$ is a ring homomorphism, and M is a left R -module. (a) Prove that M is a left S -module via $s \star x = f(s) \cdot x$ for $s \in S$ and $x \in M$ by checking the S -module axioms. (b) Give a short proof that M is a left S -module using the results of Exercise 12.
35. (a) Assume M is a left R -module annihilated by an ideal I of R . Complete the verification (from §17.7) that M is a left R/I -module via $(r + I) \bullet m = r \cdot m$ for $r \in R$ and $m \in M$. (b) Given a left R -module N and an ideal I of R , check that IN is a submodule of N .
36. Suppose R is a commutative ring with maximal ideal M . (Maximality means that for any ideal J of R with $M \subseteq J \subseteq R$, we have $J = M$ or $J = R$.) Prove that the ring R/M is a field. [Given a nonzero $x + M \in R/M$ with $x \in R$, consider the ideal $M + Rx$.]
37. Let $N = \left\{ \begin{bmatrix} 0 & b \\ 0 & d \end{bmatrix} : b, d \in \mathbb{R} \right\}$, which is an $M_2(\mathbb{R})$ -submodule of $M_2(\mathbb{R})$ by Exercise 23. Use the fundamental homomorphism theorem to prove that there is an $M_2(\mathbb{R})$ -module isomorphism $M_2(\mathbb{R})/N \cong N$.
38. Let $f : M \rightarrow N$ be an R -map between left R -modules M and N . (a) Prove: for any submodule M' of M , $f[M']$ is a submodule of N , and this submodule is contained

- in $\text{img}(f)$. (b) Prove: for any submodule N' of N , $f^{-1}[N']$ is a submodule of M , and this submodule contains $\ker(f)$.
39. Let M and N be left R -modules. Use the fundamental homomorphism theorem for modules to prove the following results. (a) $M/\{0_M\} \cong M$. (b) $M/M \cong \{0_M\}$. (c) $\frac{M \times N}{M \times \{0_N\}} \cong N$. (d) Given an index set I , left R -modules M_i for $i \in I$, and a submodule N_i of M_i for each $i \in I$, $\frac{\prod_{i \in I} M_i}{\prod_{i \in I} N_i} \cong \prod_{i \in I} (M_i/N_i)$.
40. Suppose $f : M \rightarrow N$ is a homomorphism of left R -modules, and H is a submodule of M . Attempt to define $f' : M/H \rightarrow N$ by $f'(x+H) = f(x)$ for all $x \in M$. (a) Suppose H is not contained in $\ker(f)$. Prove f' is not a well-defined function. (b) Suppose $H \subseteq \ker(f)$. Prove (as stated in §17.8) that f' is an R -module homomorphism with image $\text{img}(f)$ and kernel $\ker(f)/H$.
41. Assume H and N are submodules of a left R -module M with $H \subseteq N$. Prove: for all $x \in M$, $x+H \in N/H$ iff $x \in N$.
42. This exercise proves some results from set theory that were used in the proof of the correspondence theorem for modules. Let X and Y be sets, and let $f : X \rightarrow Y$ be any function. (a) Prove: for all $W \subseteq Y$, $f[f^{-1}[W]] \subseteq W$. (b) Give an example to show equality need not hold in (a). Then prove that equality does hold if f is surjective. (c) Prove: for all $U \subseteq X$, $f^{-1}[f[U]] \supseteq U$. (d) Give an example to show equality need not hold in (c). Then prove that equality does hold if f is one-to-one. (e) Prove: for all U_1, U_2 with $U_1 \subseteq U_2 \subseteq X$, $f[U_1] \subseteq f[U_2]$. (f) Prove: for all W_1, W_2 with $W_1 \subseteq W_2 \subseteq Y$, $f^{-1}[W_1] \subseteq f^{-1}[W_2]$.
43. Assume M is a finite left R -module. (a) Given a submodule C of M , explain why $|M| = |C| \cdot |M/C|$. (b) Let A and B be submodules of M . Use an appropriate isomorphism theorem to prove that $|A+B| = |A| \cdot |B| / |A \cap B|$.
44. We know every subgroup of \mathbb{Z} has the form $\mathbb{Z}m$ for a unique $m \geq 0$, and for all $a, b \in \mathbb{Z}$, $\mathbb{Z}a \subseteq \mathbb{Z}b$ iff b divides a . Use this information and the correspondence theorem to draw the submodule lattices of the following quotient \mathbb{Z} -modules: (a) $\mathbb{Z}/8\mathbb{Z}$; (b) $\mathbb{Z}/35\mathbb{Z}$; (c) $\mathbb{Z}/60\mathbb{Z}$. Verify by inspection of the drawings that all maximal chains of submodules have the same length.
45. Draw the lattice of all \mathbb{Z}_2 -submodules of the \mathbb{Z}_2 -module $\mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_2$.
46. Let R be the set of upper-triangular matrices in $M_2(\mathbb{Z}_2)$. (a) Check that R is a subring of $M_2(\mathbb{Z}_2)$. (b) Draw the submodule lattice of R , viewed as a left R -module. (c) Draw the submodule lattice of R , viewed as a right R -module. (d) Draw the lattice of two-sided ideals of the ring R .
47. (Solve this exercise without using the Jordan–Hölder Theorem.) (a) Suppose A is a submodule of an R -module M with the property that A and M/A are simple R -modules. Suppose B is a submodule of M different from $\{0\}$, A , and M . Prove that B and M/B are simple R -modules. (It may aid your intuition to draw a picture of what you know about the submodule lattice of M .) (b) Give an example of R, A, B, M satisfying the conditions in (a), such that A and B are non-isomorphic R -modules.
48. Fix $m, n \in \mathbb{N}^+$ with $\gcd(m, n) = 1$. Use the recognition theorem for direct products to show that the \mathbb{Z} -module $M = \mathbb{Z}_{mn}$ is isomorphic to the direct product of its submodules $\mathbb{Z}n$ and $\mathbb{Z}m$. Conclude that $\mathbb{Z}_{mn} \cong \mathbb{Z}_m \times \mathbb{Z}_n$ as \mathbb{Z} -modules.
49. *Recognition Theorem for Direct Products and Sums.* (a) Suppose N_1, N_2, \dots, N_k are submodules of a left R -module M such that $M = \sum_{i=1}^k N_i$ and

$N_i \cap (N_1 + N_2 + \cdots + N_{i-1}) = \{0\}$ for $2 \leq i \leq k$. Prove that

$g : N_1 \times N_2 \times \cdots \times N_k \rightarrow M$, given by $g(x_1, \dots, x_k) = x_1 + \cdots + x_k$ for $x_i \in N_i$,

is an R -module isomorphism. (b) Suppose I is any indexing set, and $\{N_i : i \in I\}$ is a collection of submodules of a left R -module M such that $M = \sum_{i \in I} N_i$ and $N_i \cap \sum_{j \in I, j \neq i} N_j = \{0_M\}$ for all $i \in I$. Prove that $M \cong \bigoplus_{i \in I} N_i$.

50. Find the length of the following \mathbb{Z} -modules: (a) \mathbb{Z}_{32} ; (b) \mathbb{Z}_{60} ; (c) $\mathbb{Z}_{12} \times \mathbb{Z}_{15}$; (d) \mathbb{Z} ; (e) \mathbb{Z}_p^n where p is prime and $n \geq 1$.
51. Given $n \in \mathbb{N}^+$, apply the Jordan–Hölder Theorem to the \mathbb{Z} -module \mathbb{Z}_n to prove that n can be factored into a product of prime integers that are uniquely determined by n (up to reordering).
52. Let V be a vector space over a field F . (a) Prove: if V has an n -element basis for some $n \in \mathbb{N}^+$, then there exists a maximal chain of subspaces of V of length n . (b) Prove: if V has an infinite basis, then there exist infinite chains of subspaces of V . (c) Use the Jordan–Hölder Theorem to prove that if V is finite-dimensional, then all bases of V have the same (finite) size.
53. Let (x_1, \dots, x_n) be a list of elements in a left R -module M , where R is a nonzero ring. Show that this list is R -linearly dependent in each of the following situations. (a) $x_i = 0_M$ for some i ; (b) $x_i = x_j$ for some $i \neq j$; (c) x_1, \dots, x_{n-1} generate M .
54. Let $v_1 = (2, 3)$ and $v_2 = (1, 5)$. (a) Show that (v_1, v_2) is an ordered basis of the \mathbb{R} -module $\mathbb{R} \times \mathbb{R}$. (b) Show that (v_1, v_2) is not an ordered basis of the \mathbb{Z} -module $\mathbb{Z} \times \mathbb{Z}$. (c) True or false: (v_1, v_2) is an ordered basis of the $(\mathbb{R} \times \mathbb{R})$ -module $\mathbb{R} \times \mathbb{R}$. (Explain.)
55. Let $f : N \rightarrow P$ be an R -linear map between two left R -modules N and P . (a) Prove that if f is surjective (onto) and the list (x_1, \dots, x_n) generates N , then $(f(x_1), \dots, f(x_n))$ generates P . (b) Prove that if f is injective (one-to-one) and the list (x_1, \dots, x_n) is R -linearly independent in N , then $(f(x_1), \dots, f(x_n))$ is R -linearly independent in P .
56. In the UMP for free R -modules from §17.12, check in detail that g' is an R -linear map and $g'(x) = g(x)$ for all $x \in S$.
57. Let I be an ideal of a commutative ring R such that $\{0\} \neq I \neq R$. (a) Show that the R -module R/I has no finite basis. Indicate where you use the hypotheses $\{0\} \neq I$ and $I \neq R$. (b) True or false: The R/I -module R/I has a basis. (Explain.)
58. Prove: for all $x, y \in R$, the list (x, y) is an R -basis of the left R -module R iff there exist $r, s \in R$ with $xr = 1 = ys$, $yr = 0 = xs$, and $rx + sy = 1$.
59. (a) Let A, B, C be submodules of a left R -module M such that $A \subseteq B$, $A + C = B + C$, and $A \cap C = B \cap C$. Prove that $A = B$. (b) Give an example of an \mathbb{R} -module M and submodules A, B, C such that $A + C = B + C$, $A \cap C = B \cap C$, and yet $A \neq B$.
60. Give a justified example of each of the following: (a) a ring R and an R -module M that is not finitely generated; (b) a ring R and an infinite R -module M that is finitely generated but not cyclic; (c) a ring R and a finitely generated R -module M that has no basis; (d) a list (v_1, v_2) in the \mathbb{Z} -module $\mathbb{Z} \times \mathbb{Z}$ that is \mathbb{Z} -linearly independent but does not generate $\mathbb{Z} \times \mathbb{Z}$; (e) a ring R , a commutative group $(M, +)$, and a function $\bullet : R \times M \rightarrow M$ satisfying every module axiom except $1_R \bullet x = x$ for all $x \in M$; (f) a simple left $M_3(\mathbb{Q})$ -module.

61. True or false? Explain each answer. (a) The set of all subgroups of a fixed R -module M , ordered by set inclusion, is always a complete lattice. (b) Every submodule of a product R -module $M \times N$ must have the form $A \times B$, where A is a submodule of M and B is a submodule of N . (c) For any submodules A and B of a left R -module M , there is an R -module isomorphism $(A + B)/A \cong A/(A \cap B)$. (d) A left R -module M is simple iff M has exactly two submodules. (e) There exists exactly one function $\cdot : \mathbb{Z} \times (\mathbb{Z} \times \mathbb{Z}) \rightarrow \mathbb{Z} \times \mathbb{Z}$ that turns the additive group $\mathbb{Z} \times \mathbb{Z}$ into a left \mathbb{Z} -module. (f) There exists exactly one function $\cdot : (\mathbb{Z} \times \mathbb{Z}) \times \mathbb{Z} \rightarrow \mathbb{Z}$ that turns the additive group \mathbb{Z} into a left $(\mathbb{Z} \times \mathbb{Z})$ -module. (g) Every left R -module M is a sum of cyclic submodules.
62. Suppose we are given ten left R -modules A, B, C, \dots , etc., and thirteen R -linear maps f, g, h, \dots , etc., as shown in the following diagram:

$$\begin{array}{ccccccc} A & \xrightarrow{f} & B & \xrightarrow{g} & C & \xrightarrow{h} & D & \xrightarrow{k} & E \\ \downarrow \alpha & & \downarrow \beta & & \downarrow \gamma & & \downarrow \delta & & \downarrow \epsilon \\ A' & \xrightarrow{f'} & B' & \xrightarrow{g'} & C' & \xrightarrow{h'} & D' & \xrightarrow{k'} & E' \end{array}$$

(This means $f : A \rightarrow B$, $\alpha : A \rightarrow A'$, etc.) Assume that: (1) $\text{img}(f) = \ker(g)$; (2) $\text{img}(g) = \ker(h)$; (3) $\text{img}(h) = \ker(k)$; (4) $\text{img}(f') = \ker(g')$; (5) $\text{img}(g') = \ker(h')$; (6) $\text{img}(h') = \ker(k')$; (7) $f' \circ \alpha = \beta \circ f$; (8) $g' \circ \beta = \gamma \circ g$; (9) $h' \circ \gamma = \delta \circ h$; (10) $k' \circ \delta = \epsilon \circ k$; (11) β is one-to-one; (12) δ is one-to-one; (13) α is onto. (a) Prove that γ is one-to-one. Explicitly indicate which of the 13 hypotheses are used by your proof. (b) Keep assumptions (1) through (10), but now assume: (11') β is onto; (12') δ is onto; (13') ϵ is one-to-one. Prove that γ is onto. Explicitly indicate which of the 13 hypotheses are used by your proof. (c) Deduce that, if (1) through (10) hold and α, β, δ , and ϵ are all isomorphisms, then γ is an isomorphism.

63. Let M be a left R -module, and let X be the set of simple submodules of M . For $S = \{N_i : i \in I\} \subseteq X$, define $\text{Sp}(S)$ to be the set of all simple modules $N \in X$ with $N \subseteq \sum_{i \in I} N_i$. (a) Prove that (X, Sp) is an independence structure (see Chapter 16). (b) Prove that $S = \{N_i : i \in I\}$ is independent (as defined in §16.2) iff for all $i \in I$, $N_i \cap \sum_{j \in I, j \neq i} N_j = \{0\}$ iff every $x \in \sum_{i \in I} N_i$ can be written uniquely as $x = \sum_{i \in I} x_i$ with $x_i \in N_i$ and all but finitely many x_i 's equal to zero. When this holds, we say that $\sum_{i \in I} N_i$ is the *internal direct sum* of the N_i 's, denoted $\bigoplus_{i \in I} N_i$. (c) Suppose M is the sum of all its simple submodules. Show: for $S \subseteq X$, $\text{Sp}(S) = X$ iff $\sum_{N \in S} N = M$. Deduce that $M = \bigoplus_{i \in I} N_i$ for some family of simple modules N_i . (d) Suppose $M = \bigoplus_{i \in I} N_i = \bigoplus_{j \in J} P_j$ with all N_i and P_j in X . Prove $|I| = |J|$. (e) Explain why the existence of bases for a vector space V and the uniqueness of the cardinality of bases of V are special cases of (c) and (d).

Principal Ideal Domains, Modules over PIDs, and Canonical Forms

Given a field F , we know that every finite-dimensional F -vector space V is isomorphic to F^n for some integer $n \geq 0$, where F^n is the direct product of n copies of the vector space F . Similarly, given any finitely generated commutative group G , we proved in Chapter 15 that G is isomorphic to a direct product of finitely many cyclic groups. Now, F -vector spaces are the same thing as F -modules, and commutative groups are the same thing as \mathbb{Z} -modules. So, the two theorems just mentioned provide a *classification* of all finitely generated R -modules, where R is either a field or the ring \mathbb{Z} .

In this chapter, we will prove a more general classification theorem that includes both of the previous results as special cases. To obtain this theorem, we must isolate the key properties of fields and the ring \mathbb{Z} that made the previous theorems work. This leads us to study *principal ideal domains* (or *PIDs*), which are integral domains where every ideal can be generated by a single element. We will see that \mathbb{Z} is a PID, and so is the polynomial ring $F[x]$ in one variable with coefficients in a field F . As in the case of \mathbb{Z} and $F[x]$, we will see that elements in general PIDs have unique factorizations into “irreducible” elements (which are analogous to prime integers or irreducible polynomials).

The key theorem in this chapter asserts that for any PID R , every finitely generated R -module M is isomorphic to a direct product of cyclic R -modules

$$R/Ra_1 \times R/Ra_2 \times \cdots \times R/Ra_k$$

for some $k \in \mathbb{N}$ and $a_i \in R$. As in the case of \mathbb{Z} -modules, we can arrange that the a_i 's satisfy certain additional conditions (for instance, that each a_i divide a_{i+1} , or that every nonzero a_i be a “prime power” in R). When we impose appropriate conditions of this kind, the ideals Ra_1, \dots, Ra_k appearing in the decomposition of M are *uniquely determined* by M .

The proof of this theorem will mimic the proof of the classification theorem for commutative groups, which occupies most of Chapter 15. This chapter can be read independently of that one and yields the main results of that chapter as a special case. However, the reader is urged to study that chapter first, to get used to the essential ideas of the proof in the very concrete setting of integer-valued matrices. This chapter does assume knowledge of definitions and basic facts about R -modules and free R -modules, which were covered in Chapter 17. We will also need some material on one-variable polynomials over a field from Chapter 3.

As an application of the classification theorem for general PIDs, we prove the *rational canonical form* theorem for linear operators. This theorem shows that every linear map $T : V \rightarrow V$ defined on a finite-dimensional F -vector space V can be represented (relative to an appropriate ordered basis) by a matrix with an especially simple structure. To obtain this matrix and to prove its uniqueness, we first use T to make V into a finitely generated $F[x]$ -module, and then apply the structure theorems of this chapter to this module. For fields F satisfying appropriate hypotheses, we will use similar techniques to give another

derivation of the *Jordan canonical form* theorem from Chapter 8. We also discuss Smith normal forms, rational canonical forms, and Jordan forms for matrices.

18.1 Principal Ideal Domains

Let us begin by spelling out the definition of a principal ideal domain in more detail. Recall from §1.2 that an *integral domain* is a commutative ring $(R, +, \cdot)$ such that $1_R \neq 0_R$ and R has no zero divisors other than 0. The last condition means that for all $a, b \in R$, $a \cdot b = 0_R$ implies $a = 0_R$ or $b = 0_R$. The following *cancellation law* holds in integral domains R : for all $x, y, z \in R$ with $x \neq 0_R$, if $xy = xz$ then $y = z$. To prove this, rewrite $xy = xz$ as $xy - xz = 0$ and then as $x(y - z) = 0$. Since $x \neq 0$, we get $y - z = 0$ and $y = z$.

Next, recall from §1.4 that an *ideal* of a commutative ring R is a subset I of R satisfying these closure conditions: $0_R \in I$; for all $x, y \in I$, $x + y \in I$; for all $x \in I$, $-x \in I$; and for all $x \in I$ and $r \in R$, $r \cdot x \in I$. One readily checks that for any commutative ring R and any $a \in R$, the set $Ra = \{r \cdot a : r \in R\}$ is an ideal of R containing a . This ideal is called the *principal ideal generated by a* . An ideal I of R is a *principal ideal* iff there exists $a \in R$ with $I = Ra$. A *principal ideal domain (PID)* is an integral domain R such that every ideal of R is a principal ideal.

We know that the ring \mathbb{Z} is an integral domain, since the product of any two nonzero integers is nonzero. To prove that \mathbb{Z} is a PID, consider any ideal I of \mathbb{Z} . By the first three closure conditions in the definition of an ideal, we see that I is an additive subgroup of the group $(\mathbb{Z}, +)$. In §15.1, we used integer division with remainder to prove that such a subgroup must have the form $\mathbb{Z}n = \{kn : k \in \mathbb{Z}\}$ for some integer $n \geq 0$. So I is a principal ideal.

The reader may check (Exercise 2) that every field F is a PID. Next, let us show that a *one-variable polynomial ring $F[x]$ with coefficients in a field F is a PID*. In §3.4, we used degree arguments to see that $F[x]$ is an integral domain. To show that every ideal I of $F[x]$ is principal, we use the division theorem for one-variable polynomials (§3.6). On one hand, if $I = \{0\}$, then $I = F[x]0$ is a principal ideal. On the other hand, if I is a nonzero ideal, we can choose a nonzero $g \in I$ of minimum possible degree. Because I is an ideal and $g \in I$, I contains the principal ideal $F[x]g = \{pg : p \in F[x]\}$. We now prove that $I \subseteq F[x]g$, which will imply that $I = F[x]g$ is principal. Fix any $f \in I$. Dividing f by g produces a unique quotient $q \in F[x]$ and remainder $r \in F[x]$ with $f = qg + r$ and $r = 0$ or $\deg(r) < \deg(g)$. If $r = 0$, then $f = qg \in F[x]g$ as needed. If $r \neq 0$, then $r = f + (-q)g \in I$ since $f, g \in I$ and I is an ideal. But then $\deg(r) < \deg(g)$ contradicts minimality of the degree of g , so $r \neq 0$ cannot occur.

One can check (Exercise 3) that $\mathbb{Z}[x]$ and $F[x_1, \dots, x_n]$ (where $n > 1$) are examples of integral domains that are not PIDs.

18.2 Divisibility in Commutative Rings

In any commutative ring R , we can define concepts related to divisibility by analogy with the familiar definitions for integers and polynomials. Given $a, b \in R$, we say a divides b in R and write $a|b$ iff there exists $c \in R$ with $b = ac$. When $a|b$, we also say b is a *multiple* of a and a is a *divisor* of b . One checks immediately that: for all $a \in R$, $a|a$ (reflexivity); for all

$a, b, c \in R$, if $a|b$ and $b|c$ then $a|c$ (transitivity); and for all $a \in R$, $1|a$ and $a|0$. A *unit* of the ring R is an element $u \in R$ such that $u|1$. This means that $1 = uv = vu$ for some $v \in R$, so that the units of R are precisely the invertible elements of R relative to the multiplication in R . We let R^* be the set of units of R .

The set R under the divisibility relation $|$ is almost a poset, since $|$ is reflexive and transitive on R . However, for almost all commutative rings R , antisymmetry does not hold for $|$. In other words, there can exist $x \neq y$ in R with $x|y$ and $y|x$. For all $x, y \in R$, we define x and y to be *associates* in R , denoted $x \sim y$, iff $x|y$ and $y|x$. One checks that \sim is an equivalence relation on the set R . Associate ring elements behave identically with respect to divisibility; more precisely, given $a, a', b, c \in R$ with $a \sim a'$, one checks that $a|b$ iff $a'|b$, and $c|a$ iff $c|a'$. Similarly, $u \in R$ is a unit of R iff $u \sim 1$.

In *integral domains* R , one has the following alternative description of when two elements $x, y \in R$ are associates: $x \sim y$ iff $y = ux$ for some unit $u \in R^*$. In one direction, assume $y = ux$ for some unit u of R . There is $v \in R$ with $vu = 1$, so $vy = vux = 1x = x$. Since $y = ux$ and $x = vy$, we see that $x|y$ and $y|x$ in R , hence $x \sim y$. Conversely, assume x and y are associates in R . Then $y = ax$ and $x = by$ for some $a, b \in R$. Combining these, we get $1x = by = (ba)x$ and $1y = ax = aby = (ba)y$. If either x or y is nonzero, the cancellation law for integral domains gives $ba = 1$, so that $a \in R^*$ and $y = ax$ as needed. If $x = 0 = y$, then $y = ux$ holds for the unit $u = 1_R$.

For example, in \mathbb{Z} the units are $+1$ and -1 , so $x \sim y$ in \mathbb{Z} iff $y = \pm x$. In $F[x]$ with F a field, the units are the nonzero constant polynomials, so $p \sim q$ in $F[x]$ iff $q = cp$ for some nonzero $c \in F$. It follows that each equivalence class of \sim in \mathbb{Z} contains exactly one *nonnegative* integer, whereas each equivalence class of \sim in $F[x]$ (other than $\{0\}$) contains exactly one *monic* polynomial.

Next we define common divisors, common multiples, gcds, and lcms. Let a_1, \dots, a_k be fixed elements of a commutative ring R . We say $d \in R$ is a *common divisor* of a_1, \dots, a_k iff $d|a_i$ for $1 \leq i \leq k$. We say $e \in R$ is a *common multiple* of a_1, \dots, a_k iff $a_i|e$ for $1 \leq i \leq k$. We say $d \in R$ is a *greatest common divisor (gcd)* of a_1, \dots, a_k iff d is a common divisor of a_1, \dots, a_k such that for all common divisors c of a_1, \dots, a_k , $c|d$. We say $e \in R$ is a *least common multiple (lcm)* of a_1, \dots, a_k iff e is a common multiple of a_1, \dots, a_k such that for all common multiples c of a_1, \dots, a_k , $e|c$. (Compare these definitions to the definitions of lower bounds, upper bounds, greatest lower bounds, and least upper bounds in a poset, given in the Appendix.)

We warn the reader that greatest common divisors of a_1, \dots, a_k need not exist in general. Even if gcds do exist, they usually are not unique. Indeed, given associate ring elements $d, d' \in R$, it follows from the definitions that d is a gcd of a_1, \dots, a_k iff d' is a gcd of a_1, \dots, a_k . On the other hand, if $c, d \in R$ are any two gcds of a_1, \dots, a_k , then $c \sim d$. Similar results hold for lcms. In \mathbb{Z} and $F[x]$, we get around the non-uniqueness by always using nonnegative integers and monic polynomials as our gcds and lcms.

18.3 Divisibility and Ideals

Early in the development of abstract algebra, it was realized that *divisibility of elements* in a commutative ring R can be conveniently studied by instead looking at *containment of principal ideals*. To explain this, we need the fundamental observation that *for all* $a, b \in R$, $a|b$ in R iff $Rb \subseteq Ra$. Fix $a, b \in R$. On one hand, assume $a|b$ in R , say $b = ac$ for some $c \in R$. To prove $Rb \subseteq Ra$, fix $x \in Rb$. Then $x = rb$ for some $r \in R$, hence $x = r(ac) = (rc)a \in Ra$. On the other hand, assume $Rb \subseteq Ra$. Then $b = 1b \in Rb$, so $b \in Ra$, so $b = sa$ for some

$s \in R$, so $a|b$. When using this result, one must take care to remember that the “smaller” element in the divisibility relation $a|b$ corresponds to the larger ideal in the containment relation $Rb \subseteq Ra$.

One checks readily that the set X of all ideals of R is a poset ordered by set inclusion. The subset Z of X consisting of just the *principal* ideals of R is also a poset ordered by \subseteq . Let us translate some of the definitions in the last section into statements about ideal containment. First, $a, b \in R$ are associates in R iff a and b generate the same principal ideal in R . For, $a \sim b$ iff $a|b$ and $b|a$ iff $Rb \subseteq Ra$ and $Ra \subseteq Rb$ iff $Ra = Rb$. Second, u is a unit of R iff $Ru = R$. This follows since $u \in R^*$ iff $u \sim 1$ iff $Ru = R1 = R$. Next, given $d, a_1, \dots, a_k \in R$, d is a common divisor of a_1, \dots, a_k iff $d|a_i$ for all i iff $Ra_i \subseteq Rd$ for all i iff the ideal Rd is an upper bound for the set of ideals $\{Ra_1, \dots, Ra_k\}$ in the poset Z . Similarly, the element d is a common multiple of the elements a_i iff Rd is a lower bound for $\{Ra_1, \dots, Ra_k\}$ in Z ; d is a gcd of the a_i iff Rd is the least upper bound of $\{Ra_1, \dots, Ra_k\}$ in Z ; and d is an lcm of the a_i iff Rd is the greatest lower bound of $\{Ra_1, \dots, Ra_k\}$ in Z .

We can use these remarks to give a quick proof that *for any a_1, \dots, a_k in a PID R , there exist gcds and lcms for this list of elements*. We need only find a least upper bound and a greatest lower bound for the set of ideals $S = \{Ra_1, Ra_2, \dots, Ra_k\}$ in the poset Z . Because R is a PID, Z and X are the same poset. On one hand, one sees that $I = Ra_1 + Ra_2 + \dots + Ra_k = \{r_1a_1 + r_2a_2 + \dots + r_ka_k : r_1, \dots, r_k \in R\}$ is an ideal of R that is the least upper bound for S in the poset $X = Z$. We know I is principal, so $I = Rd$ for some $d \in R$, and any generator d for I is a gcd of a_1, \dots, a_k . On the other hand, $J = Ra_1 \cap Ra_2 \cap \dots \cap Ra_k$ is an ideal of R that is the greatest lower bound for S in the poset $X = Z$. We know J is principal, so $J = Re$ for some $e \in R$, and any generator e for J is an lcm of a_1, \dots, a_k . (Recall from §17.3 that the poset of all submodules of a fixed left R -module M is a lattice. We have just reproved the special case of that result where $M = R$ is a PID.)

The proof in the previous paragraph yields a special property of gcds in PIDs: *given a_1, \dots, a_k in a PID R having gcd $d \in R$, there exist $r_1, \dots, r_k \in R$ with $d = r_1a_1 + \dots + r_ka_k$* . In other words, each gcd of a_1, \dots, a_k is an R -linear combination of the a_i 's. To prove this, recall $d \in R$ is a gcd of a_1, \dots, a_k iff Rd is the (unique) least upper bound of $\{Ra_1, \dots, Ra_k\}$ in the poset Z iff $Rd = Ra_1 + \dots + Ra_k$. So $d \in Ra_1 + \dots + Ra_k$ can be written in the required form.

18.4 Prime and Irreducible Elements

To continue our discussion of factorization theory in a commutative ring R , we need to generalize the definition of a *prime* integer or an *irreducible* polynomial. Two different generalizations are possible, but these generalizations will coincide when R is a PID. Let $p \in (R \sim R^*) \sim \{0\}$ be a nonzero element of R that is not a unit of R . We say p is a *prime* in R iff for all $f, g \in R$, $p|(fg)$ implies $p|f$ or $p|g$. We say p is *irreducible* in R iff for all $q \in R$, $q|p$ implies $q \sim p$ or $q \in R^*$. This means that the only divisors of p in R are the units of R (which divide everything) and the associates of p .

In any integral domain R , every prime p must be irreducible. To see why, let p be prime in R and assume $q \in R$ divides p . Then $p = qg$ for some $g \in R$. By primeness of p , either $p|q$ or $p|g$. If $p|q$, then (since also $q|p$) we see that $q \sim p$. On the other hand, if $p|g$, write $g = rp$ for some $r \in R$. Then $1p = qg = (qr)p$. As $p \neq 0$ in the integral domain R , we can

cancel p to get $qr = 1$, so that $q \in R^*$. On the other hand, there exist integral domains R and irreducible elements $p \in R$ that are not prime in R (see Exercise 27).

Let us translate the definitions of prime and irreducible elements into statements about principal ideals. First, the assumption that p is a nonzero non-unit means that $Rp \neq \{0\}$ and $Rp \neq R$, so that Rp is a proper nonzero ideal of R . Rewriting the definition of prime element, we see that p is prime in R iff for all $f, g \in R$, $R(fg) \subseteq Rp$ implies $Rf \subseteq Rp$ or $Rg \subseteq Rp$. Now, $R(fg) \subseteq Rp$ iff $fg \in Rp$, and similarly for Rf and Rg . So we can also say that p is prime in R iff for all $f, g \in R$, $fg \in Rp$ implies $f \in Rp$ or $g \in Rp$. In ring theory, an ideal I of a commutative ring R is called a *prime ideal* iff $I \neq R$ and for all $f, g \in R$, $fg \in I$ implies $f \in I$ or $g \in I$. So we have shown that p is a prime element in R iff Rp is a nonzero prime ideal of R .

Next, rewriting the definition of irreducible element, we get that p is irreducible in R iff for all $q \in R$, $Rp \subseteq Rq$ implies $Rq = Rp$ or $Rq = R$. (Recall that p still satisfies $\{0\} \neq Rp \neq R$.) In other words, irreducibility of a nonzero non-unit p means that the principal ideal Rp is a *maximal element* in the poset of all *proper, principal ideals* of R . In ring theory, an ideal I of a commutative ring R is called a *maximal ideal* iff $I \neq R$ and for all ideals J with $I \subseteq J$, either $I = J$ or $I = R$. In other words, maximal ideals (as just defined) are maximal elements in the poset of all *proper ideals* of R . In the case of a PID R , the poset of proper principal ideals of R is the same as the poset of proper ideals of R . We conclude that, in a PID R , p is irreducible iff Rp is a nonzero maximal ideal of R .

It is routine to prove that *every maximal ideal in any commutative ring R is a prime ideal* (Exercise 15). By the remarks in the last two paragraphs, we conclude that every irreducible element in a PID is also a prime element. Using the theorem from three paragraphs ago, we see that *irreducible elements and prime elements coincide in a PID*.

18.5 Irreducible Factorizations in PIDs

A *unique factorization domain (UFD)* is an integral domain R in which the following theorem is true: (a) For all nonzero $f \in R \sim R^*$, there exist $k \in \mathbb{N}^+$ and irreducible elements $p_1, \dots, p_k \in R$ with $f = p_1 p_2 \cdots p_k$. (b) For any two factorizations $f = u p_1 p_2 \cdots p_k = v q_1 q_2 \cdots q_m$ with all p_i and q_j irreducible in R and $u, v \in R^*$, we must have $k = m$ and, after reordering the q_j 's, $p_i \sim q_i$ for $1 \leq i \leq k$. It is well-known that \mathbb{Z} is a UFD. We proved in Chapter 3 that for every field F , $F[x]$ is a UFD. Here we will prove the celebrated theorem that *every PID is a UFD*.

To prove statement (a), we need a lemma about chains of ideals in a PID. Suppose R is a PID and we have an infinite sequence of ideals

$$I_1 \subseteq I_2 \subseteq I_3 \subseteq \cdots \subseteq I_k \subseteq \cdots.$$

Then there exists k_0 such that $I_k = I_{k_0}$ for all $k \geq k_0$. Informally, we say that *every ascending chain of ideals in a PID must stabilize*. To prove this, let $I = \bigcup_{k=1}^{\infty} I_k$ be the union of all the ideals in the sequence. It is routine to verify that I is an ideal of R , and $I_k \subseteq I$ for all $k \in \mathbb{N}^+$. Since R is a PID, there exists $a \in R$ with $I = Ra$. Now $a \in I$, so $a \in I_{k_0}$ for some fixed index k_0 . Then, for any $k \geq k_0$, $I = Ra \subseteq I_{k_0} \subseteq I_k \subseteq I$, which proves that $I = I_k = I_{k_0}$ for all such k .

We prove (a) holds for the PID R by contradiction. Assuming (a) fails, we can find a nonzero counterexample $f_0 \in R \sim R^*$ that cannot be factored into a product of irreducible elements in R . In particular, f_0 itself cannot be irreducible, so we can write $f_0 = gh$ where $g, h \in R$ are not zero, are not units of R , and are not associates of f_0 . In terms of ideals,

these conditions mean that $Rf_0 \subsetneq Rg \subsetneq R$ and $Rf_0 \subsetneq Rh \subsetneq R$. Now, if g and h could both be factored into products of irreducible elements, then f could be so factored as well by combining the two factorizations. It follows that g or h must also be a counterexample to (a). Let $f_1 = g$ if g is a counterexample, and $f_1 = h$ otherwise. Now $Rf_0 \subsetneq Rf_1 \subsetneq R$ and f_1 is a counterexample to (a). Then we can repeat the argument to produce another counterexample f_2 with $Rf_0 \subsetneq Rf_1 \subsetneq Rf_2 \subsetneq R$. This process can be continued indefinitely (using the Axiom of Choice), ultimately producing an infinite strictly ascending chain of ideals in the PID R . But this contradicts the result proved in the previous paragraph. So (a) does hold for R .

For (b), assume $f = up_1p_2 \cdots p_k = vq_1q_2 \cdots q_m$, where $u, v \in R^*$, $k, m \in \mathbb{N}^+$, and all p_i 's and q_j 's are irreducible in the PID R . Because R is a PID, all p_i 's and q_j 's are prime. In particular, the prime element p_1 divides $f = vq_1q_2 \cdots q_m$, so p_1 must divide some q_j . (Here p_1 cannot divide the unit v , or p_1 would also be a unit of R .) Reordering the q 's if needed, assume that p_1 divides q_1 . As p_1 is a non-unit and q_1 is irreducible, we obtain $p_1 \sim q_1$. Write $q_1 = wp_1$ for some $w \in R^*$; then $up_1p_2 \cdots p_k = (vw)p_1q_2 \cdots q_m$ where u and $v' = vw$ are units of R . We are in an integral domain, so we can cancel p_1 to obtain $up_2 \cdots p_k = v'q_2 \cdots q_m$. Now we repeat the argument to get $p_2 \sim q_j$ for some $j \geq 2$. We can reorder to ensure $j = 2$ and then modify the unit v' to replace q_2 by its associate p_2 . Then cancel p_2 from both sides and continue until all p_i 's have been matched with appropriate q_j 's. Note that $k < m$ is impossible, since otherwise we would obtain $u = v^*q_{k+1} \cdots q_m$ after k cancellation steps, contradicting the fact that q_m is not a unit. Similarly, $k > m$ is impossible, so $k = m$, and $p_i \sim q_i$ for $1 \leq i \leq k$ after reordering the q 's.

18.6 Free Modules over a PID

We now begin our study of finitely generated modules over a fixed PID R . As in the case of \mathbb{Z} -modules (Chapter 15), the first step is to look at properties of finitely generated *free* R -modules. (See §17.12 for proofs of the general facts about free modules recalled here.) An R -module M is *free and finitely generated* (or *f.g. free*, for short) iff M has a *finite ordered basis* $B = (v_1, \dots, v_n)$, which is an R -linearly independent list of vectors that spans the R -module M . In more detail, R -independence means that for all $c_1, \dots, c_n \in R$, if $c_1v_1 + \cdots + c_nv_n = 0_M$, then $c_1 = \cdots = c_n = 0_R$. The assertion that B spans M means that for all $w \in M$, there exist $d_1, \dots, d_n \in R$ such that $w = d_1v_1 + \cdots + d_nv_n$. The list of scalars (d_1, \dots, d_n) is uniquely determined by w , by R -independence of B ; we call (d_1, \dots, d_n) the *coordinates of w relative to the ordered basis B* . The map $T : M \rightarrow R^n$ such that $T(w) = (d_1, \dots, d_n)$ is an R -module isomorphism. Hence, every f.g. free R -module M with an n -element basis is isomorphic to the free R -module R^n whose elements are n -tuples of scalars in R .

The R -module M with ordered basis $B = (v_1, \dots, v_n)$ satisfies the following *universal mapping property* (UMP): for every R -module N and every list $w_1, \dots, w_n \in N$, there exists a unique R -linear map $U : M \rightarrow N$ with $U(v_i) = w_i$ for $1 \leq i \leq n$, namely $U(\sum_{i=1}^n d_i v_i) = \sum_{i=1}^n d_i w_i$. Using the UMP, we saw that if N is a finitely generated R -module (free or not) generated by n elements, then there is a surjective R -linear map $U : R^n \rightarrow N$, and hence an R -module isomorphism $N \cong R^n / \ker(U)$. In other words, *every finitely generated R -module is isomorphic to a quotient module of a f.g. free R -module by some submodule*. This explains why the study of f.g. free R -modules will help us understand the structure of all finitely generated R -modules.

However, as in the case of \mathbb{Z} , we will need to know that *for a PID R , any submodule P*

of any f.g. free R -module M is also f.g. free. Moreover, if M has a k -element basis, then P has a d -element basis for some $d \leq k$. We imitate the proof in §15.8. We know $M \cong R^k$ for some $k \in \mathbb{N}$, so we can assume $M = R^k$ without loss of generality. Use induction on k . If $k = 0$, then P and M must be the zero module, which is f.g. free with an empty basis. Suppose $k = 1$; the assumption that P is a submodule of R^1 means that P is an ideal of the ring R . Because R is a PID, there exists $a \in R$ with $P = Ra$. If $a = 0$, then $P = \{0\}$ is f.g. free with a basis of size zero. Otherwise, $a \neq 0$, and $B = (a)$ is a generating list for P of size 1. Is this list R -linearly independent? Given $c \in R$ with $ca = 0$, we see that $c = 0$ since $a \neq 0$ and R is an integral domain. So B is a one-element basis for P , which is therefore f.g. free. Note how the conditions in the definition of a PID were exactly what we needed to make this base case work.

Proceeding to the induction step, fix $k > 1$ and assume the theorem is known for all f.g. free R -modules having bases of size less than k . Let P be a fixed submodule of R^k . Define $Q = P \cap (R^{k-1} \times \{0\})$, which is a submodule of the free R -module $(R^{k-1} \times \{0\}) \cong R^{k-1}$. By induction, Q is f.g. free with some ordered basis (v_1, \dots, v_{d-1}) where $d-1 \leq k-1$. Consider the projection map $T : R^k \rightarrow R$ given by $T((r_1, \dots, r_k)) = r_k$ for $r_i \in R$. T is evidently R -linear, so $T[P] = \{T(x) : x \in P\}$ is an R -submodule of R . Since R is a PID, $T[P] = Ry$ for some $y \in R$. Fix an element $v_d \in P$ with $T(v_d) = y$, so v_d has last coordinate y .

If $y = 0$, then $P = Q$ is f.g. free with a basis of size $d-1 < k$. If $y \neq 0$, we will show that $B = (v_1, \dots, v_d)$ is an ordered R -basis of P , so that P is f.g. free with a basis of size $d \leq k$. First, is B an R -linearly independent list? Assume $c_1, \dots, c_d \in R$ satisfy $c_1v_1 + \dots + c_dv_d = 0$. Applying the R -linear map T and noting that $T(v_i) = 0$ for $i < d$ (since $v_1, \dots, v_{d-1} \in Q$), we get $0 = 0 + \dots + 0 + c_dT(v_d) = c_dy$. As $y \neq 0$ and R is an integral domain, $c_d = 0$ follows. Now, since $c_1v_1 + \dots + c_{d-1}v_{d-1} = 0$, the known R -linear independence of v_1, \dots, v_{d-1} gives $c_1 = \dots = c_{d-1} = 0$ as well. So B is R -linearly independent.

Second, does B span the R -module P ? Fix $z = (z_1, \dots, z_k) \in P$. Since $z_k = T(z) \in T[P] = Ry$, we have $z_k = ry$ for some $r \in R$. Then $z - rv_d$ is in the R -submodule P and has last coordinate $z_k - ry = 0$, so $z - rv_d \in Q$. Therefore $z - rv_d = e_1v_1 + \dots + e_{d-1}v_{d-1}$ for some $e_i \in R$, and we see that z itself is an R -linear combination of v_1, \dots, v_d . This completes the induction proof.

18.7 Operations on Bases

Assume R is a PID and M is a f.g. free R -module with ordered basis $X = (v_1, \dots, v_n)$. We can perform various transformations on X that produce new ordered bases for M . For example, by analogy with §15.4, there are three *elementary operations* we could apply to X . Operation (B1) interchanges v_i and v_j for some i, j ; operation (B2) replaces v_i by uv_i for some i and some unit $u \in R^*$; and operation (B3) replaces v_i by $v_i + bv_j$ for some $i \neq j$ and some $b \in R$. It can be checked that applying any finite sequence of such operations to X produces a new ordered basis of M , and each elementary operation is reversible.

We will require an even more general operation on bases that includes (B1), (B2), and (B3) as special cases. Suppose we are given a, b, c, d in the PID R such that $u = ad - bc$ is a unit of R . Operation (B4) acts on X by replacing v_i with $v'_i = av_i + bv_j$ and replacing v_j with $v'_j = cv_i + dv_j$ for some $i \neq j$ in $\{1, 2, \dots, n\}$. We can also write this as

$$\begin{bmatrix} v'_i \\ v'_j \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} v_i \\ v_j \end{bmatrix}. \quad (18.1)$$

We claim the new list $X' = (v_1, \dots, v'_i, \dots, v'_j, \dots, v_n)$ is another ordered basis for M , and X can be recovered from X' by another operation of type (B4).

By inverting the 2×2 matrix in (18.1), we get

$$\begin{bmatrix} v_i \\ v_j \end{bmatrix} = \begin{bmatrix} u^{-1}d & -u^{-1}b \\ -u^{-1}c & u^{-1}a \end{bmatrix} \begin{bmatrix} v'_i \\ v'_j \end{bmatrix}. \quad (18.2)$$

where $(u^{-1}d)(u^{-1}a) - (-u^{-1}b)(-u^{-1}c) = u^{-2}(da - bc) = u^{-1}$ is a unit in R since $u \in R^*$. This shows that we can go from X' back to X by an operation of type (B4).

Let us check that X' spans M and is linearly independent. Given $w \in M$, write $w = \sum_{k=1}^n d_k v_k$ for scalars $d_k \in R$. Using (18.2), we can replace $d_i v_i$ and $d_j v_j$ in this expression by R -linear combinations of v'_i and v'_j . So w is in the span of X' . Next, assume $0 = e_i v'_i + e_j v'_j + \sum_{k \neq i,j} e_k v_k$ for scalars $e_s \in R$. Using (18.1), this equation becomes

$$\begin{aligned} 0 &= e_i(av_i + bv_j) + e_j(cv_i + dv_j) + \sum_{k \neq i,j} e_k v_k \\ &= (e_i a + e_j c)v_i + (e_i b + e_j d)v_j + \sum_{k \neq i,j} e_k v_k. \end{aligned}$$

By the known linear independence of X , it follows that $e_k = 0$ for all $k \neq i, j$ and $e_i a + e_j c = e_i b + e_j d = 0$. In matrix terms, $[e_i \ e_j] \begin{bmatrix} a & b \\ c & d \end{bmatrix} = [0 \ 0]$. Right-multiplying by the inverse matrix gives $e_i = e_j = 0$. So X' is linearly independent.

Note that the elementary operations (B1), (B2), and (B3) are all special cases of operation (B4): for (B1), take $a = d = 0$ and $b = c = 1$; for (B2), take $a = u$, $d = 1$ (for any j), and $b = c = 0$; for (B3), take $a = d = 1$, $c = 0$ and any b . One readily checks that (B2) can still be applied to replace v_1 by uv_1 (for some $u \in R^*$) in the case $n = 1$.

18.8 Matrices of Linear Maps between Free Modules

Our next step is to study the matrices that represent linear maps between free modules. Assume R is a PID, M is a free R -module with ordered basis $X = (v_1, \dots, v_n)$, N is a free R -module with ordered basis $Y = (w_1, \dots, w_m)$, and $T : M \rightarrow N$ is a fixed R -linear map. For $1 \leq j \leq n$, we can write $T(v_j) = \sum_{i=1}^m A(i, j)w_i$ for unique scalars $A(i, j) \in R$ (which are the coordinates of $T(v_j)$ relative to the basis Y). The $m \times n$ matrix $A = (A(i, j))$ is called the *matrix of T relative to the bases X and Y* . By linearity of T , T is uniquely determined by the matrix A when X and Y are fixed and known. As in previously studied cases (when $R = \mathbb{Z}$ or R is a field), it can be shown that addition of matrices corresponds to addition of linear maps, whereas matrix multiplication corresponds to composition of linear maps.

Changing the input basis X or the output basis Y will change the matrix A that represents the given linear map T . Let us investigate how an application of basis operation (B4) to X or to Y affects the matrix A . Suppose $a, b, c, d \in R$ satisfy $u = ad - bc \in R^*$. Define $a' = u^{-1}d$, $b' = -u^{-1}b$, $c' = -u^{-1}c$, and $d' = u^{-1}a$ as in (18.2). Recall that $A^{[k]}$ denotes the k 'th column of A , whereas $A_{[k]}$ denotes the k 'th row of A . We make two claims.

- Suppose X' is obtained from X by replacing v_i and v_j by $v'_i = av_i + bv_j$ and $v'_j = cv_i + dv_j$. Then the matrix B of T relative to the bases X' and Y satisfies

$B^{[i]} = aA^{[i]} + bA^{[j]}$, $B^{[j]} = cA^{[i]} + dA^{[j]}$, and $B^{[k]} = A^{[k]}$ for all $k \neq i, j$. (We say B is obtained from A by a *type 4 column operation* on columns i and j .)

2. Suppose Y' is obtained from Y by replacing w_i and w_j by $w'_i = aw_i + bw_j$ and $w'_j = cw_i + dw_j$. Then the matrix C of T relative to the bases X and Y' satisfies $C_{[i]} = A_{[i]}a' + A_{[j]}c'$, $C_{[j]} = A_{[i]}b' + A_{[j]}d'$, and $C_{[k]} = A_{[k]}$ for all $k \neq i, j$. (We say C is obtained from A by a *type 4 row operation* on rows i and j .)

We prove the second claim, leaving the first claim as Exercise 32. For fixed $p \in \{1, \dots, n\}$, we must find the unique scalars $C(s, p) \in R$ satisfying

$$T(v_p) = C(i, p)w'_i + C(j, p)w'_j + \sum_{k \neq i, j} C(k, p)w_k.$$

We know

$$T(v_p) = A(i, p)w_i + A(j, p)w_j + \sum_{k \neq i, j} A(k, p)w_k.$$

Recalling from (18.2) that $w_i = a'w'_i + b'w'_j$ and $w_j = c'w'_i + d'w'_j$, we get

$$T(v_p) = (A(i, p)a' + A(j, p)c')w'_i + (A(i, p)b' + A(j, p)d')w'_j + \sum_{k \neq i, j} A(k, p)w_k.$$

Comparing coefficients, we see that $C(i, p) = A(i, p)a' + A(j, p)c'$, $C(j, p) = A(i, p)b' + A(j, p)d'$, and $C(k, p) = A(k, p)$ for all $k \neq i, j$. This holds for all p , so the rows of C are related to the rows of A as stated in the claim.

As special cases of the result for (B4), we can deduce how elementary column and row operations on the matrix A correspond to elementary operations of types (B1), (B2), and (B3) on the input basis X and the output basis Y . Specifically, one checks that:

3. Switching columns i and j in A corresponds to switching v_i and v_j in the input basis X ; whereas switching rows i and j of A corresponds to switching w_i and w_j in the output basis Y .
4. Multiplying column i of A by a unit $u \in R^*$ corresponds to replacing v_i by uv_i in the input basis X ; whereas multiplying row i of A by $u \in R^*$ corresponds to replacing w_i by $u^{-1}w_i$ in the output basis Y .
5. Adding b times column j of A to column i of A (where $b \in R$) corresponds to replacing v_i by $v_i + bv_j$ in the input basis X ; whereas adding b times row j of A to row i of A corresponds to replacing w_j by $w_j - bw_i$ in the output basis Y .

We will also need the following properties, which the reader is asked to prove in Exercise 32.

6. Suppose $e \in R$ divides every entry of a matrix $A \in M_{m,n}(R)$. If we apply any sequence of type 4 row and column operations to A , then e will still divide every entry of the new matrix.
7. Suppose B is obtained from A by applying a type 4 column operation as in claim 1. Then $B = AV$, where $V \in M_n(R)$ is an invertible matrix with entries $V(i, i) = a$, $V(j, i) = b$, $V(i, j) = c$, $V(j, j) = d$, $V(k, k) = 1_R$ for all $k \neq i, j$, and all other entries of V are zero (cf. §4.7 and §4.9).
8. Suppose C is obtained from A by applying a type 4 row operation as in claim 2. Then $C = UA$, where $U \in M_m(R)$ is an invertible matrix with entries $U(i, i) = a'$, $U(i, j) = c'$, $U(j, i) = b'$, $U(j, j) = d'$, $U(k, k) = 1_R$ for all $k \neq i, j$, and all other entries of U are zero (cf. §4.8 and §4.9).

18.9 Reduction Theorem for Matrices over a PID

We now have all the tools we need to prove the following matrix reduction theorem. Let A be an $m \times n$ matrix with entries in a PID R . There exists a finite sequence of type 4 row and column operations on A that will reduce A to a new matrix

$$B = \begin{bmatrix} a_1 & 0 & 0 & \dots & 0 \\ 0 & a_2 & 0 & \dots & 0 \\ 0 & 0 & a_3 & \dots & 0 \\ \vdots & & & \ddots & \end{bmatrix}, \quad (18.3)$$

in which there are $s \geq 0$ nonzero elements $a_1, \dots, a_s \in R$ on the main diagonal, a_i divides a_{i+1} in R for $1 \leq i < s$, and all other entries of B are zero. For brevity, we write $B = \text{diag}(a_1, \dots, a_s)_{m \times n}$, omitting the $m \times n$ when $s = m = n$. In terms of ideals, the divisibility conditions on the a_j 's are equivalent to $Ra_1 \supseteq Ra_2 \supseteq Ra_3 \supseteq \dots \supseteq Ra_s$. Later, we will prove that s and the ideals satisfying this containment condition are uniquely determined by A , and hence the a_j 's are unique up to associates in R . So the matrix B is essentially unique; it is called a *Smith normal form* of A . By repeated use of properties 7 and 8 from §18.8, we see that $B = PAQ$ for some invertible $P \in M_m(R)$ and $Q \in M_n(R)$, where P (resp. Q) is the product of all the matrices U (resp. V) used to accomplish the type 4 row (resp. column) operations needed to reduce A to B .

We proved the reduction theorem for $R = \mathbb{Z}$ in §15.11. If we try to repeat that proof in the setting of general PIDs, a problem emerges. The old proof made heavy use of integer division with remainder, as well as the fact that there is no infinite strictly decreasing sequence of positive integers. These proof ingredients can be generalized to a class of rings called *Euclidean domains* (defined in Exercise 5), but they are not available in all PIDs. To execute the proof at this level of generality, a new trick is needed.

Recall that all PIDs are UFDs, so that every nonzero non-unit $a \in R$ has a factorization $a = p_1 p_2 \cdots p_s$ into irreducible elements p_i in R ; moreover, any other irreducible factorization $a = q_1 q_2 \cdots q_t$ has $s = t$ and $p_i \sim q_i$ after appropriate reordering. Define the *length of a in R* to be $\text{len}(a) = s$, the number of factors appearing in any irreducible factorization of a . For any unit u of R , let $\text{len}(u) = 0$; $\text{len}(0_R)$ is undefined. For any nonzero matrix A with entries in R , let $\text{len}(A) \in \mathbb{N}$ be the minimum length of all the nonzero entries of A . Given nonzero $a, b, d \in R$, one may check that: $\text{len}(ab) = \text{len}(a) + \text{len}(b)$; if d divides a in R , then $\text{len}(d) \leq \text{len}(a)$ with equality iff $d \sim a$; and if $d = \gcd(a, b)$ where a does not divide b , then $\text{len}(d) < \text{len}(a)$.

We now begin the proof of the matrix reduction theorem when R is a PID. Observe that the theorem certainly holds if $A = 0$ or $m = 0$ or $n = 0$, so we assume $m, n > 0$ and $A \neq 0$ throughout the rest of the proof. Using induction on m (the number of rows), we can assume the theorem is known to hold for all matrices with fewer than m rows.

Step 1: We show we can apply finitely many type 4 row and column operations to A to produce a matrix A_1 such that some nonzero entry in A_1 divides all entries of A_1 . The proof uses induction on $\text{len}(A)$. If $\text{len}(A) = 0$, then some entry of A is a unit of R , which divides all elements of R and hence divides all entries of A . Next, assume $\text{len}(A) = \ell > 0$ and the result of Step 1 is known to hold for all matrices of length less than ℓ . Let $e = A(i, j)$ be a nonzero entry of A with $\text{len}(e) = \ell$. If e happens to divide all entries of A , then the conclusion of Step 1 already holds for the matrix A .

Suppose instead that there exists at least one entry of A not divisible by $e = A(i, j)$. Case 1: There is $k \neq j$ such that e does not divide $f = A(i, k)$. Since R is a PID, we know $g = \gcd(e, f)$ exists in R , and $g = ae + bf$ for some $a, b \in R$. We have $e = gd$ and $f = gc$ for

some $c, d \in R$. Cancelling $g \neq 0$ in $g = a(gd) + b(gc)$ gives $1_R = ad + bc$. So, we can apply a type 4 column operation to A that replaces $A^{[j]}$ by $aA^{[j]} + bA^{[k]}$ and $A^{[k]}$ by $-cA^{[j]} + dA^{[k]}$. The new matrix A' has i, j -entry $ae + bf = g$, and $\text{len}(g) < \text{len}(e)$ since e does not divide f . So $\text{len}(A') \leq \text{len}(g) < \text{len}(A)$. By induction, we can apply further reduction steps to A' to achieve the conclusion of Step 1.

Case 2: There is $k \neq i$ such that $e = A(i, j)$ does not divide $f = A(k, j)$. We argue as in Case 1, but this time we use a type 4 row operation to replace e by $g = \gcd(e, f)$, which lowers the length of the matrix.

Case 3: $e = A(i, j)$ divides everything in row i and column j , but for some $i_1 \neq i$ and $j_1 \neq j$, e does not divide $f = A(i_1, j_1)$. Pictorially, rows i, i_1 and columns j, j_1 look like this for some $u, v \in R$:

$$\begin{bmatrix} e & \cdots & ue \\ \vdots & & \vdots \\ ve & \cdots & f \end{bmatrix}.$$

Adding $(1 - v)$ times row i to row i_1 produces:

$$\begin{bmatrix} e & \cdots & ue \\ \vdots & & \vdots \\ e & \cdots & f + (1 - v)ue \end{bmatrix}.$$

If this new matrix has lower length than A , we are done by induction. Otherwise, note e cannot divide $f + (1 - v)ue$ in R , lest e divide f . So Case 1 now applies to row i_1 , and we can complete Step 1 as in that case.

Step 2: Let $a_1 = A_1(i, j)$ be a nonzero entry in A_1 dividing all entries of A_1 . We show A_1 can be further reduced to the form

$$A_2 = \begin{bmatrix} a_1 & 0 & \cdots & 0 \\ 0 & & & \\ \vdots & & A' & \\ 0 & & & \end{bmatrix}, \quad (18.4)$$

where $A' \in M_{m-1, n-1}(R)$, and a_1 divides all entries of A' in R . To prove this, recall that applying type 4 row and column operations to A_1 will never change the property that a_1 divides every entry of the matrix. To begin, bring a_1 into the 1, 1-position by switching row 1 and row i and switching column 1 and column j . Since a_1 divides every entry in row 1, we can subtract appropriate multiples of column 1 from each later column to make the other entries in row 1 become zero. Similarly, we can use row operations to produce zeroes below a_1 in the first column. The current matrix now looks like (18.4). Since a_1 still divides all entries of the full matrix, a_1 divides every entry of A' .

Step 3: We show how A_2 can be reduced to the normal form (18.3). Since A' has $m - 1 < m$ rows, the induction hypothesis shows that we can apply type 4 row and column operations to this $(m - 1) \times (n - 1)$ matrix to obtain a matrix in normal form with entries a_2, \dots, a_s on the main diagonal, all other entries zero, and $a_i | a_{i+1}$ in R for $2 \leq i < s$. We can apply the same type 4 operations to the full matrix A_2 , and these operations will not disturb the zeroes we have already created in row 1 and column 1. Furthermore, a_1 will continue to divide all entries of the matrix throughout the reduction of A_2 . In particular, at the end, a_1 will divide a_2 , and we have reached the required normal form for A .

18.10 Structure Theorems for Linear Maps and Modules

The matrix reduction theorem proved above translates into the following *structure theorem for linear maps between free R -modules*. Let R be a PID, let N and M be f.g. free R -modules, and let $T : N \rightarrow M$ be an R -linear map. There exist an ordered basis $X = (x_1, \dots, x_n)$ for N and an ordered basis $Y = (y_1, \dots, y_m)$ for M such that the matrix of T relative to X and Y is in Smith normal form (18.3). So, there exist $s \in \mathbb{N}$ and nonzero $a_1, \dots, a_s \in R$ with $a_i | a_{i+1}$ for $1 \leq i < s$, $T(x_i) = a_i y_i$ for $1 \leq i \leq s$, and $T(x_i) = 0_M$ for $s < i \leq n$. (Later, we will prove that s and the ideals $Ra_1 \supseteq Ra_2 \supseteq \dots \supseteq Ra_s \neq \{0\}$ are uniquely determined by T .)

To prove existence, start with any ordered bases X_0 for N and Y_0 for M , and let A be the matrix of T relative to X_0 and Y_0 . Use a finite sequence of type 4 operations to bring A into Smith normal form. To ensure that each new matrix still represents T , we perform the appropriate basis operation (B4) on the input basis (when we do a column operation on A) or on the output basis (when we do a row operation on A), using the rules explained in §18.8. At the end, we will have new ordered bases X for N and Y for M satisfying the required properties.

Next we prove the existence part of the fundamental structure theorem for modules over a PID: *for any finitely generated module M over a PID R , there exist $d \in \mathbb{N}$ and $b_1, \dots, b_d \in R$ with $R \neq Rb_1 \supseteq Rb_2 \supseteq \dots \supseteq Rb_d \supseteq \{0\}$ and*

$$M \cong R/Rb_1 \times R/Rb_2 \times \dots \times R/Rb_d. \quad (18.5)$$

(Note that some b_j 's might be zero here, in which case $R/Rb_j \cong R$.) Each ideal Rb_j is called an *invariant factor* of M ; the generators b_j of these ideals are also called invariant factors. In §18.14, we will prove that d and the sequence of ideals $(Rb_1, Rb_2, \dots, Rb_d)$ are uniquely determined by M .

To prove the existence part of the theorem, suppose M is an R -module generated by m elements. Recall (§18.6) that $M \cong R^m/P$ for some submodule P of R^m , where P is a free R -module with a basis of size $n \leq m$. The inclusion map $T : P \rightarrow R^m$ given by $T(x) = x$ for all $x \in P$ is an R -linear map between free R -modules. So there exist an ordered basis $X = (x_1, \dots, x_n)$ of P , an ordered basis $Y = (y_1, \dots, y_m)$ for R^m , $s \in \mathbb{N}$, and nonzero $a_1, \dots, a_s \in R$ with $Ra_1 \supseteq \dots \supseteq Ra_s \neq \{0\}$, $T(x_i) = a_i y_i$ for $1 \leq i \leq s$, and $T(x_i) = 0$ for $s < i \leq n$. Now $T(x_i) = x_i \neq 0$ for all i (since X is a basis), so we must have $s = n$ and $x_i = a_i y_i$ for $1 \leq i \leq n$. Define $a_i = 0$ for $n < i \leq m$, so $Ra_1 \supseteq \dots \supseteq Ra_m \supseteq \{0\}$. As in §15.13, the UMP for free R -modules provides an isomorphism $R^m \cong R^m$ sending y_i to e_i (the standard basis vector) for $1 \leq i \leq m$, and this isomorphism sends P to $P_1 = Ra_1 \times Ra_2 \times \dots \times Ra_m$. So

$$M \cong R^m/P \cong R^m/P_1 \cong (R/Ra_1) \times (R/Ra_2) \times \dots \times (R/Ra_m),$$

where the last step uses the fundamental homomorphism theorem for modules (see Chapter 17, Exercise 39(d)). To finish, we need only delete any initial factors R/Ra_i that are equal to zero (which happens iff $Ra_i = R$ iff a_i is a unit of R).

Continuing to imitate §15.13, we now derive a “prime power” version of the structure theorem for modules. We need this lemma: *suppose R is a PID and $a \in R$ has irreducible factorization $a = p_1^{e_1} \cdots p_k^{e_k}$ where p_1, \dots, p_k are non-associate irreducible elements of R and $e_1, \dots, e_k \in \mathbb{N}^+$; then there is an R -module isomorphism*

$$R/Ra \cong (R/Rp_1^{e_1}) \times \dots \times (R/Rp_k^{e_k}).$$

To prove this, define $T : R \rightarrow \prod_{i=1}^k (R/Rp_i^{e_i})$ by $T(x) = (x + Rp_1^{e_1}, \dots, x + Rp_k^{e_k})$ for $x \in R$. The map T is R -linear, and $x \in \ker(T)$ iff $x + Rp_i^{e_i} = 0$ for all i iff $p_i^{e_i}|x$ for all i iff $a = \text{lcm}(p_1^{e_1}, \dots, p_k^{e_k})|x$ iff $x \in Ra$. So T induces an isomorphism $T' : R/Ra \rightarrow \text{img}(T)$. It now suffices to show T is onto, which will be accomplished by showing that each generator $(0, \dots, 1 + Rp_i^{e_i}, \dots, 0)$ of $\prod_{i=1}^k (R/Rp_i^{e_i})$ is in the image of T . Note $r = p_i^{e_i}$ and $s = \prod_{j \neq i} p_j^{e_j}$ have $\text{gcd } 1_R$ (by comparing unique prime factorizations), so there exist $b, c \in R$ with $br + cs = 1_R$. Consider $T(cs)$. For any $k \neq i$, the coset $cs + Rp_k^{e_k}$ is zero since $p_k^{e_k}$ divides s . On the other hand, the coset $cs + Rp_i^{e_i} = (1 - br) + Rp_i^{e_i} = 1 + Rp_i^{e_i}$ since $-br + Rp_i^{e_i}$ is the zero coset. Thus, $T(cs) = (0, \dots, 1 + Rp_i^{e_i}, \dots, 0)$ as needed.

Applying the lemma to each nonzero b_j in (18.5), we obtain the following structural result: *for any finitely generated module M over a PID R , there exist $k \in \mathbb{N}$ and $q_1, \dots, q_k \in R$ such that each q_i is either zero or $p_i^{e_i}$ for some irreducible $p_i \in R$ and $e_i \in \mathbb{N}^+$, and*

$$M \cong R/Rq_1 \times R/Rq_2 \times \cdots \times R/Rq_k. \quad (18.6)$$

The ideals Rq_j (as well as the generators q_j of these ideals) are called the *elementary divisors* of the R -module M . In §18.15, we will prove that these ideals (counted with multiplicity) are uniquely determined by the module M .

18.11 Minors and Matrix Invariants

We want to prove the uniqueness of the Smith normal form of a matrix $A \in M_{m,n}(R)$ or a linear map between f.g. free R -modules. Before doing so, we need some preliminary results on *matrix invariants*, which are quantities depending on A that do not change when we multiply A on the left or right by an invertible matrix with entries in R .

Let A be an $m \times n$ matrix with entries in a PID R . We will need to consider submatrices of A formed by keeping only certain rows and columns of A . More precisely, given a subset $I = \{i_1 < i_2 < \cdots < i_k\}$ of $[m] = \{1, 2, \dots, m\}$ and a subset $J = \{j_1 < j_2 < \cdots < j_\ell\}$ of $[n] = \{1, 2, \dots, n\}$, the *submatrix of A with rows in I and columns in J* is the matrix $A_{I,J}$ with r, s -entry $A_{I,J}(r, s) = A(i_r, j_s)$ for $1 \leq r \leq k$ and $1 \leq s \leq \ell$. For example, $A_{[m], \{2,3,5\}}$ is the submatrix obtained by keeping all m rows of A and columns 2, 3, and 5 of A . Using this notation, the Cauchy–Binet Formula (proved in §5.14) can be stated as follows: Given $U \in M_{k,m}(R)$ and $V \in M_{m,k}(R)$ with $k \leq m$,

$$\det(UV) = \sum_{L \subseteq [m], |L|=k} \det(U_{[k],L}) \det(V_{L,[k]}).$$

Fix k with $1 \leq k \leq \min(m, n)$. A $k \times k$ *submatrix of A* is a matrix $A_{I,J}$ with $|I| = |J| = k$. For each choice of I and J of size k , let $d_{I,J} = \det(A_{I,J}) \in R$; $d_{I,J}$ is called the k 'th order minor of A indexed by I and J . Let $g_k(A)$ be a gcd in R of all the determinants $d_{I,J}$ with $|I| = |J| = k$. We claim that for all invertible $P \in M_m(R)$ and all invertible $Q \in M_n(R)$, $g_k(A) \sim g_k(PAQ)$ in R . In terms of ideals, this says that $Rg_k(A) = Rg_k(PAQ)$, so that the ideal generated by any gcd of all the k 'th order minors of A is a matrix invariant of A .

First, we use the Cauchy–Binet formula to prove that $g_k(A) \sim g_k(PA)$. Fix $I = \{i_1 < \cdots < i_k\} \subseteq [m]$ and $J = \{j_1 < \cdots < j_k\} \subseteq [n]$ of size k . Note that $(PA)_{I,J} = P_{I,[m]} A_{[m],J}$, since the r, s -entry of both sides is $(PA)(i_r, j_s) = \sum_{t=1}^m P(i_r, t)A(t, j_s)$. Since $k \leq m$, we can apply the Cauchy–Binet formula to the $k \times m$ matrix $U = P_{I,[m]}$ and the

$m \times k$ matrix $V = A_{[m],J}$. We obtain

$$\det((PA)_{I,J}) = \det(UV) = \sum_{\substack{L \subseteq [m], \\ |L|=k}} \det(U_{[k],L}) \det(V_{L,[k]}) = \sum_{\substack{L \subseteq [m], \\ |L|=k}} \det(P_{I,L}) \det(A_{L,J}).$$

This formula shows that the minor $\det((PA)_{I,J})$ of the matrix PA is an R -linear combination of various k 'th order minors $d_{L,J}$ of the matrix A . Therefore, if $e \in R$ is a common divisor of all the k 'th order minors of A , then e divides each k 'th order minor of PA . In particular, $e = g_k(A)$ divides all minors $\det(A_{I,J})$, so $g_k(A)$ is a common divisor of all minors $\det((PA)_{I,J})$, so $g_k(A)$ divides $g_k(PA) = \gcd\{\det((PA)_{I,J}) : |I| = |J| = k\}$. We have now proved $g_k(A)|g_k(PA)$ for all $A \in M_{m,n}(R)$ and all invertible $P \in M_m(R)$. Applying this result with A replaced by PA and P replaced by P^{-1} , we see that $g_k(PA)|g_k(P^{-1}(PA)) = g_k(A)$. Hence, $g_k(A) \sim g_k(PA)$ as needed.

By a similar argument (Exercise 48), we can use the Cauchy–Binet formula to show $g_k(A) \sim g_k(AQ)$ for any $A \in M_{m,n}(R)$ and any invertible $Q \in M_n(R)$. So $g_k(A) \sim g_k(PAQ)$ when P and Q are invertible over R . (One can also prove $g_k(PA) \sim g_k(A) \sim g_k(AQ)$ without appealing to the Cauchy–Binet formula, but instead invoking multilinearity properties of determinants — see Exercise 49.)

18.12 Uniqueness of Smith Normal Form

Let R be a PID, and let A be an $m \times n$ matrix with entries in R . In §18.9, we proved that there exist invertible matrices $P \in M_m(R)$ and $Q \in M_n(R)$ such that $B = PAQ$ has the form

$$B = \text{diag}(a_1, a_2, \dots, a_s)_{m \times n},$$

where $s \in \mathbb{N}$, $a_1, \dots, a_s \in R$, and $Ra_1 \supseteq Ra_2 \supseteq \dots \supseteq Ra_s \neq \{0\}$.

Our goal here is to prove the *uniqueness* of s and the ideals Ra_1, \dots, Ra_s satisfying the properties just stated. More specifically, we will show that for any invertible $P' \in M_m(R)$ and $Q' \in M_n(R)$ such that

$$B' = P' A Q' = \text{diag}(b_1, b_2, \dots, b_t)_{m \times n}$$

for some $t \in \mathbb{N}$ and $b_1, \dots, b_t \in R$ with $Rb_1 \supseteq Rb_2 \supseteq \dots \supseteq Rb_t \neq \{0\}$, we must have $s = t$ and $Ra_i = Rb_i$ for $1 \leq i \leq s$ (equivalently, a_i and b_i are associates in R for all i). The ideals Ra_i (as well as their generators a_i) are called the *invariant factors* of A , and s is called the *rank* of A .

Define $a_i = 0$ for $s < i \leq \min(m, n)$ and $b_j = 0$ for $t < j \leq \min(m, n)$. Fix k with $1 \leq k \leq \min(m, n)$. On one hand, we have seen that

$$g_k(B) = g_k(PAQ) \sim g_k(A) \sim g_k(P'AQ') = g_k(B'). \quad (18.7)$$

On the other hand, we can use the special form of B to compute $g_k(B)$ directly from the definition. If $k > s$, every $k \times k$ submatrix $B_{I,J}$ must have a row and column of zeroes, so every k 'th order minor $\det(B_{I,J})$ is zero. Then $g_k(B) = 0$ since this is the gcd of a list of zeroes. Now suppose $k \leq s$. One readily sees that every k 'th order minor $\det(B_{I,J})$ is either zero or is some product of the form $a_{i_1}a_{i_2}\cdots a_{i_k} \neq 0$, where $1 \leq i_1 < i_2 < \dots < i_k \leq s$. Furthermore, one of these minors is $\det(B_{[k],[k]}) = a_1a_2\cdots a_k \neq 0$. Since $a_i|a_j$ for all $1 \leq i \leq j \leq s$, we see that $a_1a_2\cdots a_k$ divides all the k 'th order minors of B . So this ring

element is a gcd of all of these minors, and we can therefore take $g_k(B) = a_1 a_2 \cdots a_k \neq 0$ for $1 \leq k \leq s$. Letting $g_0(B) = 1_R$, we see that $a_k = g_k(B)/g_{k-1}(B)$ for $1 \leq k \leq s$. (More precisely, a_k is the unique x in the integral domain R solving $g_{k-1}(B)x = g_k(B)$.) Replacing $g_k(B)$ or $g_{k-1}(B)$ by associate ring elements will replace a_k by an associate of a_k in R .

Applying the same reasoning to B' , we see that $g_k(B') = 0$ for all $k > t$, $g_k(B') = b_1 b_2 \cdots b_k \neq 0$ for $0 \leq k \leq t$, and $b_k = g_k(B')/g_{k-1}(B')$ for $1 \leq k \leq t$. Returning to (18.7), we now see that $s = t = \max\{k : g_k(A) \neq 0_R\}$ and $a_k \sim g_k(A)/g_{k-1}(A) \sim b_k$ for $1 \leq k \leq s$, which completes the uniqueness proof.

We can deduce a similar uniqueness result for R -linear maps between f.g. free R -modules. Given a PID R , f.g. free R -modules N and M , and an R -linear map $T : N \rightarrow M$, there exist unique $s \in \mathbb{N}$ and ideals $Ra_1 \supseteq Ra_2 \supseteq \cdots \supseteq Ra_s \neq \{0\}$ such that for some ordered bases $X = (x_1, \dots, x_n)$ for N and $Y = (y_1, \dots, y_m)$ for M , $T(x_i) = a_i y_i$ for $1 \leq i \leq s$ and $T(x_i) = 0_M$ for $s < i \leq n$. Existence of one choice of s, a_1, \dots, a_s, X, Y was shown earlier (§18.10). For uniqueness, suppose $s', Ra'_1, \dots, Ra'_{s'}, X', Y'$ also satisfied all of the conclusions above. Let A be the matrix of T relative to the bases X and Y ; then A is in Smith normal form (18.3). Similarly, the matrix A' of T relative to X' and Y' is in Smith normal form with the elements a'_j on its diagonal. Considering transition matrices between the bases X and X' and the bases Y and Y' , one sees that $A' = PAQ$ for some invertible $P \in M_m(R)$ and $Q \in M_n(R)$ (Exercise 35). Applying the uniqueness result for matrices proved above, we obtain $s = s'$ and $a_i \sim a'_i$ (hence $Ra_i = Ra'_i$) for $1 \leq i \leq s$.

18.13 Torsion Submodules

Our next task is to prove the uniqueness of the ideals Rb_i and Rq_j appearing in the decompositions (18.5) and (18.6). We begin in this section by showing how to split off the “free part” of a finitely generated module over a PID.

Given an integral domain R and any R -module M , the *torsion submodule* of M is

$$\text{tor}(M) = \{x \in M : \text{for some } r \in R, r \neq 0 \text{ and } r \cdot x = 0\}.$$

To see that $\text{tor}(M)$ really is a submodule, first note $1_R \neq 0_R$ and $1_R \cdot 0_M = 0_M$, so $0_M \in \text{tor}(M)$. Next, fix $x, y \in \text{tor}(M)$ and $t \in R$. Choose nonzero $r, s \in R$ with $rx = 0 = sy$. Then $rs \neq 0$ since R is an integral domain, and $(rs) \cdot (x + y) = (rs) \cdot x + (rs) \cdot y = s \cdot (r \cdot x) + r \cdot (s \cdot y) = s0 + r0 = 0$, so $x + y \in \text{tor}(M)$. Also $r \cdot (t \cdot x) = t \cdot (r \cdot x) = t \cdot 0 = 0$ since R is commutative, so $t \cdot x \in \text{tor}(M)$.

Now suppose M and N are R -modules and $f : M \rightarrow N$ is an R -module isomorphism. One readily checks that $f[\text{tor}(M)] = \text{tor}(N)$, so that f restricts to an isomorphism from $\text{tor}(M)$ to $\text{tor}(N)$. It follows from this and the fundamental homomorphism theorem that f induces a module isomorphism $f' : M/\text{tor}(M) \rightarrow N/\text{tor}(N)$ given by $f'(x + \text{tor}(M)) = f(x) + \text{tor}(N)$ for $x \in M$. To summarize: if $M \cong N$, then $\text{tor}(M) \cong \text{tor}(N)$ and $M/\text{tor}(M) \cong N/\text{tor}(N)$.

For example, consider an R -module

$$P = R/Ra_1 \times R/Ra_2 \times \cdots \times R/Ra_k \times R^d,$$

where a_1, \dots, a_k are nonzero elements of R and $d \in \mathbb{N}$. We claim

$$\text{tor}(P) = R/Ra_1 \times \cdots \times R/Ra_k \times \{0_{R^d}\}.$$

A typical element of the right side is $z = (x_1 + Ra_1, \dots, x_k + Ra_k, 0)$ where each $x_i \in R$.

Let $r = a_1 a_2 \cdots a_k \neq 0$; note rx_i is divisible by a_i , so $r(x_i + Ra_i) = rx_i + Ra_i = 0 + Ra_i$ for $1 \leq i \leq k$, so that $rz = (0, \dots, 0, 0)$ and hence $z \in \text{tor}(P)$. On the other hand, consider $y = (x_1 + Ra_1, \dots, x_k + Ra_k, (r_1, \dots, r_d)) \in P$ with some $r_j \neq 0_R$. Multiplying y by any nonzero $s \in R$ will produce $sy = (sx_1 + Ra_1, \dots, sx_k + Ra_k, (sr_1, \dots, sr_d))$, where $sr_j \neq 0_R$. So sy cannot be 0, hence $y \notin \text{tor}(P)$. This proves the claim. From the claim, we readily deduce that $P/\text{tor}(P) \cong R^d$.

The preceding remarks immediately yield the following **splitting theorem** that lets us break apart the “free piece” and the “torsion piece” of a finitely generated module. *Suppose R is a PID and M is a finitely generated R -module such that*

$$R/Ra_1 \times \cdots \times R/Ra_k \times R^d \cong M \cong R/Rb_1 \times \cdots \times R/Rb_\ell \times R^e,$$

where every Ra_i and Rb_j is a nonzero ideal of R and $d, e \in \mathbb{N}$. Then

$$R/Ra_1 \times \cdots \times R/Ra_k \cong \text{tor}(M) \cong R/Rb_1 \times \cdots \times R/Rb_\ell \text{ and } R^d \cong M/\text{tor}(M) \cong R^e,$$

and hence $d = e$. To see why $d = e$ follows, note $M/\text{tor}(M)$ is a free R -module having a basis of size d (since this module is isomorphic to R^d) and a basis of size e (since this module is isomorphic to R^e). Since the PID R is commutative, $d = e$ follows from the theorem proved in §17.13. We call d the *Betti number* of M .

18.14 Uniqueness of Invariant Factors

We are now ready to prove the uniqueness of the sequence of ideals (the invariant factors) appearing in the decomposition (18.5). Using the splitting theorem just proved to separate out all factors of the form $R/R0 \cong R$, it will suffice to prove the following statement: *suppose R is a PID and $a_1, \dots, a_k, b_1, \dots, b_\ell$ satisfy*

$$\begin{aligned} R \neq Ra_k \supseteq \cdots \supseteq Ra_1 \neq \{0\}, \quad R \neq Rb_\ell \supseteq \cdots \supseteq Rb_1 \neq \{0\}, \\ \text{and } R/Ra_k \times \cdots \times R/Ra_1 \cong R/Rb_\ell \times \cdots \times R/Rb_1. \end{aligned} \tag{18.8}$$

Then $k = \ell$ and $Ra_i = Rb_i$ for $1 \leq i \leq k$. (We have reversed the indexing order of the a_i 's and b_j 's to simplify notation in the induction proof below.)

All quotient modules appearing here are nonzero (as $Ra_i \neq R \neq Rb_j$), so $k = 0$ iff $\ell = 0$. Assume $k, \ell > 0$. We first prove that $Ra_1 = Rb_1$ using the following ideas. Given any R -module M , the *annihilator* of M is

$$\text{ann}_R(M) = \{r \in R : \text{for all } x \in M, rx = 0\}.$$

One checks that for any commutative ring R , $\text{ann}_R(M)$ is an ideal of R . Moreover, $M \cong M'$ implies $\text{ann}_R(M) = \text{ann}_R(M')$, so that *isomorphic* R -modules have *equal* annihilators.

Given $M = R/Ra_k \times \cdots \times R/Ra_1$ as above, let us show that $\text{ann}_R(M) = Ra_1$. A typical element of M is a k -tuple of cosets $x = (x_k + Ra_k, \dots, x_1 + Ra_1)$ with all $x_i \in R$. Multiplying x by $ra_1 \in Ra_1$ (where $r \in R$) produces $(ra_1)x = (ra_1x_k + Ra_k, \dots, ra_1x_1 + Ra_1)$. Every ra_1x_i is in Ra_1 , which is contained in all the other ideals Ra_i by assumption. So $ra_1x_i + Ra_i = 0 + Ra_i$ for $1 \leq i \leq k$, proving that $(ra_1)x = 0$. This means that $Ra_1 \subseteq \text{ann}_R(M)$. For the reverse inclusion, fix $s \in \text{ann}_R(M)$. Then $s \cdot (0, \dots, 0, 1 + Ra_1) = 0_M$, so that $s + Ra_1 = 0 + Ra_1$, so that $s \in Ra_1$. The same reasoning shows that $\text{ann}_R(R/Rb_\ell \times \cdots \times R/Rb_1)$ is Rb_1 . By the result in the last paragraph, $Ra_1 = Rb_1$ follows.

Fix i with $1 \leq i - 1 \leq \min(k, \ell)$, and make the induction hypothesis that $Ra_1 = Rb_1, Ra_2 = Rb_2, \dots, Ra_{i-1} = Rb_{i-1}$. We now prove that $k \geq i$ iff $\ell \geq i$, in which case $Ra_i = Rb_i$. Assume $k \geq i$; we will show $\ell \geq i$ and $b_i|a_i$ in R . The proof will require facts about the length of a module proved in §17.11. Let $P = R/Ra_{i-1} \times \cdots \times R/Ra_1$, which appears in both of the product modules (18.8) thanks to the induction hypothesis. One may verify that $\text{len}(R/Ra) = \text{len}(a)$ for any nonzero a in a PID R (Exercise 52), so that $\text{len}(P) = \text{len}(a_1) + \cdots + \text{len}(a_{i-1}) < \infty$. Now, if we had $\ell = i - 1$, then (18.8) says

$$[R/Ra_k \times \cdots \times R/Ra_i] \times P \cong P,$$

where the term in brackets is a nonzero module Q . On one hand, the isomorphic modules $Q \times P$ and P have the same finite length. On the other hand, $\text{len}(Q \times P) > \text{len}(P)$ since $Q \neq \{0\}$. This contradiction shows $\ell \geq i$.

Note that for any R -module N and any $c \in R$ (where R is a commutative ring), $cN = \{c \cdot n : n \in N\}$ is a submodule of N ; and if $N \cong N'$ are isomorphic R -modules, then $cN \cong cN'$. Furthermore, for a direct product $N = N_1 \times N_2 \times \cdots \times N_k$, we have $cN = (cN_1) \times (cN_2) \times \cdots \times (cN_k)$. Taking $c = a_i$ and applying these remarks to (18.8), we get an isomorphism

$$[a_i(R/Ra_k) \times \cdots \times a_i(R/Ra_i)] \times a_iP \cong [a_i(R/Rb_\ell) \times \cdots \times a_i(R/Rb_i)] \times a_iP. \quad (18.9)$$

We know a_i lies in all the ideals $Ra_i \subseteq Ra_{i+1} \subseteq \cdots \subseteq Ra_k$. It follows that the product in brackets on the left side of (18.9) is the zero module. Comparing lengths of both sides (noting that $\text{len}(a_iP) \leq \text{len}(P) < \infty$), we conclude that the product in brackets on the right side must also be the zero module. In particular, $a_i(R/Rb_i) = \{0\}$, so $a_i \cdot (1 + Rb_i) = 0 + Rb_i$. This means $a_i \in Rb_i$, so $b_i|a_i$ in R .

Now, by interchanging the roles of the two product modules in (18.8), we prove similarly that $\ell \geq i$ implies $k \geq i$ and $a_i|b_i$ in R . So $k \geq i$ iff $\ell \geq i$, in which case $b_i|a_i$ and $a_i|b_i$, hence $a_i \sim b_i$, hence $Ra_i = Rb_i$. This completes the induction step. Taking $i = \min(k, \ell) + 1$, we see that $k = \ell$ and $Ra_j = Rb_j$ for $1 \leq j \leq k$.

18.15 Uniqueness of Elementary Divisors

We want to prove the uniqueness (up to reordering) of the elementary divisors appearing in the decomposition (18.6). By invoking the splitting theorem from §18.13 to remove all factors of the form $R/R0 \cong R$, it suffices to prove the following statement: *suppose R is a PID and $q_1, \dots, q_k, r_1, \dots, r_\ell$ are positive powers of irreducible elements in R such that*

$$R/Rq_1 \times \cdots \times R/Rq_k \cong R/Rr_1 \times \cdots \times R/Rr_\ell. \quad (18.10)$$

Then $k = \ell$ and the list of ideals (Rq_1, \dots, Rq_k) is a rearrangement of the list (Rr_1, \dots, Rr_ℓ) .

To simplify the proof, let $M = [Rq_1, \dots, Rq_k]$ be the set of all rearrangements of the list (Rq_1, \dots, Rq_k) ; we call M a *multiset* of ideals. This word indicates that the order in which we list the ideals is unimportant, but the number of times each ideal occurs is significant. Let X be the set of all such multisets arising from lists of finitely many ideals Rq_j with each $q_j = p_j^{e_j}$ for some irreducible $p_j \in R$ and $e_j \in \mathbb{N}^+$. Let Y be the set of all finite lists of ideals (Ra_1, \dots, Ra_m) with $R \neq Ra_1 \supseteq Ra_2 \supseteq \cdots \supseteq Ra_m \neq \{0\}$. The idea of the proof is to define bijections $f : X \rightarrow Y$ and $g : Y \rightarrow X$ that will let us invoke the known uniqueness of the invariant factors of a module. We saw this idea in the simpler setting of commutative groups in §15.19.

Let Z be a fixed set of irreducible elements in R such that no two elements of Z are associates, but every irreducible element in R is associate to some element of Z . The map g acts on $L = (Ra_1, \dots, Ra_m) \in Y$ as follows. We know each a_i factors in R into a product $u \prod_{j=1}^{n_i} p_{ij}^{e_{ij}}$ where $u \in R^*$, $p_{i1}, p_{i2}, \dots, p_{in_i}$ are distinct irreducible elements in Z and $e_{ij} \in \mathbb{N}^+$. Define $g(L)$ to be the multiset in X consisting of all ideals $Rp_{ij}^{e_{ij}}$ for $1 \leq i \leq m$ and $1 \leq j \leq n_i$. Using the lemma from §18.10, note that $\prod_{i=1}^m R/Ra_i \cong \prod_{i=1}^m \prod_{j=1}^{n_i} R/Rp_{ij}^{e_{ij}}$, no matter what order we list the terms in the direct product on the right side.

The map f acts on $M = [Rq_1, \dots, Rq_k] \in X$ as follows. Let $p_1, \dots, p_n \in Z$ be the distinct irreducible elements such that each $q_i \sim p_j^{e_j}$ for some $1 \leq j \leq n$. Place the elements of M in a matrix such that row j contains all the ideals Rq_i with $q_i \sim p_j^{e_j}$, listed with multiplicities so that the exponents e_j weakly increase reading from left to right. Suppose the longest row in the matrix has length m . Pad all shorter rows with copies of 1_R on the left so that all rows have length m . Define $f(M) = (Ra_1, Ra_2, \dots, Ra_m)$, where a_k is the product of all q_i 's appearing in column k . One may check that $f(M) \in Y$, since the construction ensures that $a_i | a_{i+1}$ for all $i < m$. Since splitting the a_k 's back into prime powers will reproduce the q_i 's in some order, we see that $g(f(M)) = M$ for all $M \in X$, and $\prod_{i=1}^k R/Rq_i \cong \prod_{j=1}^m R/Ra_j$. (One can also check that $f(g(L)) = L$ for all $L \in Y$, though we do not need this fact below.)

To begin the uniqueness proof, assume we have an isomorphism as in (18.10). Let $M_1 = [Rq_1, \dots, Rq_k]$, $M_2 = [Rr_1, \dots, Rr_\ell]$, $L_1 = f(M_1) = [Ra_1, \dots, Ra_m]$, and $L_2 = f(M_2) = [Rb_1, \dots, Rb_n]$. We have seen that $\prod_{j=1}^m R/Ra_j \cong \prod_{i=1}^k R/Rq_i \cong \prod_{i=1}^\ell R/Rr_i \cong \prod_{j=1}^n R/Rb_j$. Since $L_1, L_2 \in Y$, the uniqueness result for invariant factors shows that $L_1 = L_2$. Then $M_1 = g(f(M_1)) = g(L_1) = g(L_2) = g(f(M_2)) = M_2$, which proves the required uniqueness result for elementary divisors.

18.16 $F[x]$ -Module Defined by a Linear Operator

In the rest of this chapter, we apply the structure theorems for finitely generated modules over PIDs to derive results on canonical forms of matrices and linear operators on a vector space. Throughout, we let F be a field and V be an n -dimensional vector space over F . We also fix an F -linear map $T : V \rightarrow V$. We will define a class of matrices in $M_n(F)$ called *rational canonical forms* and show that each T is represented (relative to an appropriate ordered basis of V) by exactly one of these matrices.

To obtain this result from the preceding theory, we will use T to turn the vector space (F -module) V into an $F[x]$ -module. The addition in the $F[x]$ -module V is the given addition in the vector space V . For $v \in V$ and $p = \sum_{i=0}^d p_i x^i \in F[x]$, define scalar multiplication by $p \cdot v = \sum_{i=0}^d p_i T^i(v)$, where $T^0 = \text{id}_V$ and T^i denotes the composition of i copies of T .

Let us check the $F[x]$ -module axioms. The five additive axioms are already known to hold. For $p \in F[x]$ and $v \in V$ as above, $p \cdot v$ is in V , since T and each T^i map V to V and V is closed under addition and multiplication by scalars in F . Next, $1_{F[x]} \cdot v = 1T^0(v) = v$. Given $q = \sum_{i \geq 0} q_i x^i \in F[x]$, note $qp = \sum_{i \geq 0} \left(\sum_{k=0}^i q_k p_{i-k} \right) x^i$. Using linearity of T and

its powers, we compute:

$$\begin{aligned} (qp) \cdot v &= \sum_{i \geq 0} \left(\sum_{k=0}^i q_k p_{i-k} \right) T^i(v) = \sum_{i \geq 0} \sum_{k=0}^i q_k T^k(p_{i-k} T^{i-k}(v)) = \sum_{k \geq 0} \sum_{j \geq 0} q_k T^k(p_j T^j(v)) \\ &= \sum_{k \geq 0} q_k T^k \left(\sum_{j \geq 0} p_j T^j(v) \right) = q \cdot \left(\sum_{j \geq 0} p_j T^j(v) \right) = q \cdot (p \cdot v). \end{aligned}$$

Next,

$$(p+q) \cdot v = \sum_{i \geq 0} (p_i + q_i) T^i(v) = \sum_{i \geq 0} p_i T^i(v) + \sum_{i \geq 0} q_i T^i(v) = p \cdot v + q \cdot v.$$

Finally, given $w \in V$,

$$p \cdot (v+w) = \sum_{i \geq 0} p_i T^i(v+w) = \sum_{i \geq 0} p_i [T^i(v) + T^i(w)] = \sum_{i \geq 0} p_i T^i(v) + \sum_{i \geq 0} p_i T^i(w) = p \cdot v + p \cdot w.$$

Let us spell out some definitions from module theory in the setting of the particular $F[x]$ -module V determined by the linear operator T . First, what is an $F[x]$ -submodule of V ? This is an additive subgroup W of V such that $p \cdot w \in W$ for all $p \in F[x]$ and all $w \in W$. Taking p to be a constant polynomial, we see that a submodule W must be closed under multiplication by scalars in F . Taking $p = x$, we see that a submodule W must satisfy $x \cdot w = T(w) \in W$ for all $w \in W$. Conversely, suppose W is a subspace such that $T[W] \subseteq W$; a subspace satisfying this condition is called a *T -invariant* subspace of V . By induction on i , we see that $T^i(w) \in W$ for all $w \in W$ and all $i \geq 0$. Since W is a subspace, we then see that $p \cdot w = \sum_{i \geq 0} p_i T^i(w) \in W$ for all $w \in W$ and all $p \in F[x]$. So, *submodules of the $F[x]$ -module V are the same thing as T -invariant subspaces of V* .

Second, is V *finitely generated* as an $F[x]$ -module? We know V is finitely generated as an F -module, since the vector space V has an n -element basis $B = \{v_1, \dots, v_n\}$. We claim B also generates the $F[x]$ -module V . For, given $v \in V$, write $v = c_1 v_1 + \dots + c_n v_n$ for some $c_1, \dots, c_n \in F$. Each c_i is also a constant polynomial in $F[x]$, so v has been expressed as an $F[x]$ -linear combination of the v_i 's. (V could very well be generated, as an $F[x]$ -module, by a proper subset of B , since we could also act on each v_i by non-constant polynomials.)

Third, what does a *cyclic* $F[x]$ -submodule of V look like? Recall this is a submodule of the form $W = F[x]z$ for some fixed $z \in W$; z is called a *generator* of the submodule. Such a submodule is also called a *T -cyclic subspace* of V . Define a map $g : F[x] \rightarrow W$ by $g(p) = p \cdot z$ for all $p \in F[x]$. One may check that g is a surjective $F[x]$ -module homomorphism, which induces an $F[x]$ -module isomorphism $g' : F[x]/\ker(g) \rightarrow W$. The kernel of g is a submodule (ideal) of $F[x]$, called the *T -annihilator* of z . Since $F[x]$ is a PID, $\ker(g) = F[x]h$ for some $h \in F[x]$. Now h cannot be zero, since otherwise W would be isomorphic to $F[x]$ as an $F[x]$ -module, hence also isomorphic to $F[x]$ as an F -module. But $F[x]$ is an infinite-dimensional F -vector space and W is finite-dimensional. So $h \neq 0$, and we can take h to be the unique *monic* generator of $\ker(g)$. Write $h = h_0 + h_1 x + \dots + h_d x^d$, where $d \in \mathbb{N}$, each $h_i \in F$, and $h_d = 1$.

So far, we know that $F[x]/F[x]h$ and $W = F[x]z$ are isomorphic (both as $F[x]$ -modules and F -modules) via the map g' sending $p + F[x]h$ to $p \cdot z$ for all $p \in F[x]$. Now, using polynomial division with remainder (cf. §3.19), one checks that

$$(1 + F[x]h, x + F[x]h, x^2 + F[x]h, \dots, x^{d-1} + F[x]h)$$

is an ordered basis for the F -vector space $F[x]/F[x]h$. Applying the F -isomorphism g' to

this basis, we conclude that

$$B_z = (1 \cdot z, x \cdot z, x^2 \cdot z, \dots, x^{d-1} \cdot z) = (z, T(z), T^2(z), \dots, T^{d-1}(z))$$

is an ordered F -basis for the subspace W of V .

Since W is T -invariant, we know T restricts to a linear map $T|_W : W \rightarrow W$. What is the matrix of $T|_W$ relative to the ordered basis B_z ? Note $T|_W(z) = T(z)$, which has coordinates $(0, 1, 0, \dots, 0)$ relative to B_z . Note $T|_W(T(z)) = T(T(z)) = T^2(z)$, which has coordinates $(0, 0, 1, 0, \dots, 0)$ relative to the basis B_z . Similarly, $T|_W(T^j(z)) = T^{j+1}(z)$ for $0 \leq j < d-1$. But, when we apply $T|_W$ to the final element in B_z , $T|_W(T^{d-1}(z)) = T^d(z)$ is not in the basis. As $h \in \ker(g)$, we know $0 = h \cdot z = \sum_{i=0}^{d-1} h_i T^i(z) + T^d(z)$, so the coordinates of $T^d(z)$ relative to B_z must be $(-h_0, -h_1, \dots, -h_{d-1})$. In conclusion, the matrix we want is

$$[T|_W]_{B_z} = \begin{bmatrix} 0 & 0 & 0 & \dots & 0 & -h_0 \\ 1 & 0 & 0 & \dots & 0 & -h_1 \\ 0 & 1 & 0 & \dots & 0 & -h_2 \\ 0 & 0 & 1 & \dots & 0 & -h_3 \\ \dots & \dots & \dots & & \dots & \dots \\ 0 & 0 & 0 & \dots & 1 & -h_{d-1} \end{bmatrix}_{d \times d}. \quad (18.11)$$

This matrix is called the *companion matrix* of the monic polynomial h and is denoted C_h . Conversely, let W be any T -invariant subspace such that for some $z \in W$, $B_z = (z, T(z), \dots, T^{d-1}(z))$ is an F -basis of W and $[T|_W]_{B_z} = C_h$; it then follows that $W = F[x]z$ is a T -cyclic subspace isomorphic to $F[x]/F[x]h$ (Exercise 63(a)).

18.17 Rational Canonical Form of a Linear Map

As in the last section, let V be an n -dimensional vector space over a field F , and let $T : V \rightarrow V$ be a fixed linear map. Make V into an $F[x]$ -module via T , as described above. We know $F[x]$ is a PID and V is a finitely generated $F[x]$ -module, so the fundamental structure theorem for modules over a PID (see (18.5)) gives us an $F[x]$ -module isomorphism

$$\phi : V \rightarrow F[x]/F[x]h_1 \times F[x]/F[x]h_2 \times \dots \times F[x]/F[x]h_k \quad (h_1, \dots, h_k \in F[x]) \quad (18.12)$$

for uniquely determined ideals

$$F[x] \neq F[x]h_1 \supseteq F[x]h_2 \supseteq \dots \supseteq F[x]h_k \supseteq \{0\}.$$

Since V is finite-dimensional as an F -module, none of the ideals $F[x]h_j$ can be zero. So we can pick unique monic generators h_1, \dots, h_k of these ideals with respective degrees $d_1, \dots, d_k > 0$. These polynomials satisfy $h_i|h_{i+1}$ in $F[x]$ for $1 \leq i < k$ and are called the *invariant factors* of T .

Call the product module on the right side of (18.12) V' . For $1 \leq i \leq k$, let W'_i be the submodule $\{0\} \times \dots \times F[x]/F[x]h_i \times \dots \times \{0\}$ of V' . Note W'_i is a cyclic $F[x]$ -module generated by $z'_i = (0, \dots, 1 + F[x]h_i, \dots, 0)$. Letting $W_i = \phi^{-1}[W'_i]$ and $z_i = \phi^{-1}(z'_i)$ for each i , we obtain cyclic submodules $W_i = F[x]z_i$ of V such that $W_i \cong W'_i \cong F[x]/F[x]h_i$. Since h_i is evidently the monic polynomial of least degree sending the coset $1 + F[x]h_i$ to $0 + F[x]h_i$, we see that $F[x]h_i$ is the T -annihilator of z_i for all i .

It is routine to check that we can get an ordered basis for the F -vector space V' by concatenating ordered bases for the submodules W'_i corresponding to each factor in

the direct product. Applying ϕ^{-1} , we get an ordered basis B for the F -vector space V by concatenating ordered bases for W_1, \dots, W_k . Using the bases B_{z_1}, \dots, B_{z_k} constructed in §18.16, we get an ordered basis

$$B = (z_1, T(z_1), \dots, T^{d_1-1}(z_1), z_2, T(z_2), \dots, T^{d_2-1}(z_2), \dots, T^{d_k-1}(z_k))$$

for V . The matrix of T relative to the basis B is the block-diagonal matrix

$$[T]_B = \begin{bmatrix} C_{h_1} & 0 & \dots & 0 \\ 0 & C_{h_2} & \dots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \dots & C_{h_k} \end{bmatrix}, \quad (18.13)$$

where C_{h_i} is the companion matrix of h_i .

In general, given any square matrices A_1, \dots, A_k , we write $\text{blk-diag}(A_1, \dots, A_k)$ for the block-diagonal matrix with diagonal blocks A_1, \dots, A_k . A *rational canonical form* is a matrix $A \in M_n(F)$ of the form $A = \text{blk-diag}(C_{h_1}, \dots, C_{h_k})$, where: $k \in \mathbb{N}^+$, $h_1, \dots, h_k \in F[x]$ are monic, non-constant polynomials, and $h_i | h_{i+1}$ in $F[x]$ for $1 \leq i < k$. The matrix in (18.13) is called the *rational canonical form of T* , and the polynomials h_i are called the *invariant factors of T* .

By reversing the steps used to go from V' to $[T]_B$, we can show that *every linear map $T : V \rightarrow V$ is represented by exactly one matrix in rational canonical form*. For suppose there were another ordered basis B^* of V such that $[T]_{B^*} = \text{blk-diag}(C_{g_1}, \dots, C_{g_\ell})$, where $g_1, \dots, g_\ell \in F[x]$ are monic polynomials in $F[x]$ with $g_i | g_{i+1}$ for $1 \leq i < \ell$ and $\deg(g_i) = d_i^* > 0$. Let W_1^* be the subspace of V generated by the first d_1^* vectors in B^* , W_2^* the subspace generated by the next d_2^* vectors in B^* , and so on. The form of the matrix $[T]_{B^*}$ shows that each W_i^* is a cyclic $F[x]$ -submodule of V annihilated by g_i , namely $W_i^* = F[x]z_i^* \cong F[x]/F[x]g_i$, where z_i^* is the first basis vector in the generating list for W_i^* . One may now check (Exercise 63) that

$$V \cong W_1^* \times \cdots \times W_\ell^* \cong F[x]/F[x]g_1 \times \cdots \times F[x]/F[x]g_\ell$$

as F -modules and $F[x]$ -modules. The uniqueness result proved in §18.14 now gives $\ell = k$ and $F[x]h_i = F[x]g_i$ for all i , hence $h_i = g_i$ for all i since h_i and g_i are monic.

18.18 Jordan Canonical Form of a Linear Map

We still assume $T : V \rightarrow V$ is a fixed linear map on an n -dimensional F -vector space V . We know (see (18.6)) that there is an $F[x]$ -module isomorphism

$$\psi : V \cong V'' = F[x]/F[x]q_1 \times \cdots \times F[x]/F[x]q_s$$

where each $q_i = p_i^{e_i}$ for some monic irreducible $p_i \in F[x]$ and some $e_i \in \mathbb{N}^+$; we saw in §18.15 that s and the multiset of q_i 's are uniquely determined by these conditions. The q_i 's are called *elementary divisors* of the linear map T .

By exactly the same argument used to derive (18.13) from (18.12), we see that V is the direct sum of T -cyclic subspaces W_1, \dots, W_s with $W_i = F[x]y_i \cong F[x]/F[x]q_i$ for some $y_i \in V$, and V has an ordered basis B such that $[T]_B = \text{blk-diag}(C_{q_1}, \dots, C_{q_s})$. To obtain further structural results, we now impose the additional hypothesis that *every p_i has degree*

1, say $p_i = x - c_i$ where $c_i \in F$. For example, this hypothesis will automatically hold when $F = \mathbb{C}$, or when F is any algebraically closed field (which means all irreducible polynomials in $F[x]$ have degree 1).

Our goal is to change the basis of each T -cyclic subspace $W_i = F[x]y_i$ to obtain an even nicer matrix than the companion matrix $C_{q_i} = C_{(x-c_i)^{e_i}}$. Fix $c \in F$ and $e \in \mathbb{N}^+$, and consider any T -cyclic subspace $W = F[x]y$ of V with ordered basis $B_y = (y, T(y), \dots, T^{e-1}(y))$, such that $[T|_W]_{B_y} = C_{(x-c)^e}$. For $0 \leq k \leq e$, define $y_k = (x - c)^k \cdot y \in W$. Since $(x - c)^k = (x - c)(x - c)^{k-1}$, we have $y_0 = y$, $y_e = 0$, and for $1 \leq k \leq e$,

$$y_k = (x - c) \cdot y_{k-1} = T(y_{k-1}) - cy_{k-1}, \text{ hence } T(y_{k-1}) = y_k + cy_{k-1}.$$

Using these facts, a routine induction argument shows that the list (y_0, \dots, y_k) spans the same F -subspace as the list $(y, T(y), \dots, T^k(y))$ for $0 \leq k < e$, so $B = (y_{e-1}, \dots, y_2, y_1, y_0)$ is another ordered F -basis for W . Let us compute the matrix $[T|_W]_B$. First, $T(y_{e-1}) = y_e + cy_{e-1} = cy_{e-1}$, which has coordinates $(c, 0, \dots, 0)$ relative to B . Next, for $1 < j \leq e$, $T(y_{e-j}) = 1y_{e-j+1} + cy_{e-j}$, so that column j of the matrix has a 1 in row $j-1$, a c in row j , and zeroes elsewhere. In other words, $[T|_W]_B$ is the *Jordan block*

$$J(c; e) = \begin{bmatrix} c & 1 & 0 & \cdots & 0 \\ 0 & c & 1 & \cdots & 0 \\ 0 & 0 & c & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 1 \\ 0 & 0 & 0 & \cdots & c \end{bmatrix}.$$

(In the case $e = 1$, this is a 1×1 matrix with sole entry c .)

By concatenating bases of the form (y_{e-1}, \dots, y_0) for all the T -cyclic subspaces W_i , we obtain an ordered basis for V such that the matrix of T relative to this basis is

$$\mathbf{J} = \begin{bmatrix} J(c_1; e_1) & 0 & \cdots & 0 \\ 0 & J(c_2; e_2) & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & J(c_s; e_s) \end{bmatrix}. \quad (18.14)$$

Any block-diagonal matrix in $M_n(F)$ with Jordan blocks on the diagonal is called a *Jordan canonical form*; the particular matrix \mathbf{J} is called a *Jordan canonical form of the linear map T* . We have just proved that *every linear map T is represented by at least one matrix in Jordan canonical form, assuming that every elementary divisor of T has the form $(x - c)^e$ for some $c \in F$ and $e \in \mathbb{N}^+$ (which always holds for $F = \mathbb{C}$)*.

The Jordan canonical form for T may not be unique, since we can obtain other Jordan canonical forms by reordering the multiset of elementary divisors of T , which leads to a permutation of the Jordan blocks of \mathbf{J} . But, arguing as in §18.17, one can reverse the steps used to pass from V' to \mathbf{J} to prove that *all Jordan canonical forms of T are obtained from \mathbf{J} by reordering the Jordan blocks*. The only new detail is showing that for any subspace W of V such that $[T|_W]_B = J(c; e)$ for some basis B of W , W is a T -cyclic subspace with $W = F[x]y \cong F[x]/F[x](x - c)^e$ for some $y \in W$ (Exercise 65).

18.19 Canonical Forms of Matrices

Let $A \in M_n(F)$ be an $n \times n$ matrix with entries in a field F . We use A to define a linear map $T : F^n \rightarrow F^n$ such that $T(v) = Av$ for all $v \in F^n$. Then F^n becomes an $F[x]$ -module

via T , where the action is given by $p \cdot v = \sum_{i=0}^d p_i A^i v$ for $p = \sum_{i=0}^d p_i x^i \in F[x]$ and $v \in F^n$. The matrix of T relative to the standard ordered basis of F^n is A . We know (Chapter 6) that the matrix C of T relative to any other ordered basis of F^n has the form $C = S^{-1}AS$ for some invertible $S \in M_n(F)$; in other words, C is *similar* to A . Conversely, every matrix similar to A is the matrix of T relative to some ordered basis of F^n .

Since T has exactly one rational canonical form, it follows that *every matrix $A \in M_n(F)$ is similar to exactly one matrix of the form blk-diag(C_{h_1}, \dots, C_{h_k}), where $h_1, \dots, h_k \in F[x]$ are monic non-constant polynomials with $h_i | h_{i+1}$ for $1 \leq i < k$.* This matrix is called the *rational canonical form of A* , and the polynomials h_1, \dots, h_k are the *invariant factors of A* . Similarly, for algebraically closed fields F , *every matrix $A \in M_n(F)$ is similar to a matrix of the form blk-diag($J(c_1; e_1), \dots, J(c_s; e_s)$), and the only matrices of this form that are similar to A are those obtained by reordering the Jordan blocks.* These matrices are called *Jordan canonical forms of A* . Two matrices $A, B \in M_n(F)$ are similar iff they have the same rational canonical form iff they have the same Jordan canonical forms.

Recall that the *minimal polynomial* of $A \in M_n(F)$ is the unique monic polynomial $m_A \in F[x]$ of minimum degree such that $m_A(A) = 0$, and the *characteristic polynomial* of A is $\chi_A = \det(xI_n - A) \in F[x]$. One checks readily that similar matrices have the same minimal polynomial and the same characteristic polynomial. For monic $h \in F[x]$, one can show directly from these definitions that $m_{C_h} = \chi_{C_h} = h$ (Exercise 66). For any block-diagonal matrix $B = \text{blk-diag}(B_1, \dots, B_k)$, one sees that $m_B = \text{lcm}(m_{B_1}, \dots, m_{B_k})$ and $\chi_B = \prod_{i=1}^k \chi_{B_i}$. In the special case where B is the rational canonical form of A , we see that $\chi_A = \chi_B = \prod_{i=1}^k h_i$ and $m_A = m_B = h_k$. We have just reproved the *Cayley–Hamilton theorem*, which states that $m_A | \chi_A$ in $F[x]$, or equivalently $\chi_A(A) = 0$.

Our final theorem of the chapter states that *the invariant factors of positive degree in the Smith normal form of the matrix $xI_n - A \in M_n(F[x])$ are exactly the invariant factors in the rational canonical form of the matrix $A \in M_n(F)$.* Because of this theorem, we can use the matrix reduction algorithm in §18.9 (applied to the matrix $xI_n - A$) to calculate the rational canonical form of a given matrix A . Alternatively, we can invoke the results in §18.12 to give specific formulas for each invariant factor (involving quotients of gcds of appropriate minors of $xI_n - A$) and to show once again that the invariant factors in the rational canonical form of A are uniquely determined by A .

To prove the theorem, first consider the case where $A = C_h$ is the companion matrix of a monic polynomial $h \in F[x]$ of degree n . Using elementary row and column operations, the matrix $xI_n - C_h \in M_n(F[x])$ can be reduced to the diagonal matrix $\text{diag}(1, \dots, 1, h)_{n \times n}$ (Exercise 66). This matrix is a Smith normal form of $xI_n - C_h$. For a general matrix A , we know there is an invertible matrix $S \in M_n(F)$ such that $S^{-1}AS$ is the unique rational canonical form of A . Suppose $S^{-1}AS = \text{blk-diag}(C_{h_1}, \dots, C_{h_k})$, where h_1, \dots, h_k are the invariant factors of A . Write $\deg(h_i) = d_i > 0$ for $1 \leq i \leq k$. It follows that $S^{-1}(xI_n - A)S = xI_n - S^{-1}AS = \text{blk-diag}(xI_{d_1} - C_{h_1}, \dots, xI_{d_k} - C_{h_k})$. By the observation above, for each block $xI_{d_i} - C_{h_i}$ there are invertible $P_i, Q_i \in M_{d_i}(F[x])$ such that $P_i(xI_{d_i} - C_{h_i})Q_i = \text{diag}(1, \dots, 1, h_i)_{d_i \times d_i}$. Letting $P = \text{blk-diag}(P_1, \dots, P_k)$ and $Q = \text{blk-diag}(Q_1, \dots, Q_k)$, $PS^{-1}(xI_n - A)SQ$ will be a diagonal matrix where the diagonal entries (in some order) are h_1, \dots, h_k , and $n - k$ ones. By performing row and column interchanges, we can arrange that the diagonal consists of all the ones followed by h_1, \dots, h_k . So we have found invertible $P', Q' \in M_n(F[x])$ such that $P'(xI_n - A)Q' = \text{diag}(1, \dots, 1, h_1, \dots, h_k)_{n \times n}$. This matrix is a Smith normal form of $xI_n - A$. We see that the non-constant invariant factors in this normal form are exactly h_1, \dots, h_k , which are the invariant factors appearing in the rational canonical form of A .

18.20 Summary

1. *Divisibility Definitions.* Given x, y in a commutative ring R : x divides y (written $x|y$) iff $y = rx$ for some $r \in R$ iff $Ry \subseteq Rx$; x is an *associate* of y (written $x \sim y$) iff $x|y$ and $y|x$ iff $Rx = Ry$; x is a *unit* of R (written $x \in R^*$) iff $xy = 1_R$ for some $y \in R$ iff $x \sim 1$; x is a *zero divisor* of R iff $x \neq 0$ and there is $y \neq 0$ in R with $xy = 0_R$. When R is an integral domain, $x \sim y$ iff there exists $u \in R^*$ with $y = ux$. For $a_1, \dots, a_k, d \in R$, d is a *gcd* of a_1, \dots, a_k iff d divides all a_i , and every common divisor of the a_i 's divides d ; d is an *lcm* of a_1, \dots, a_k iff all a_i divide d , and d divides every common multiple of the a_i .
2. *Types of Ideals.* An ideal I in a commutative ring R is *principal* iff there exists $c \in R$ with $I = Rc = \{rc : r \in R\}$. An ideal I of R is *prime* iff $I \neq R$ and for all $x, y \in R$, $xy \in I$ implies $x \in I$ or $y \in I$. An ideal I of R is *maximal* iff $I \neq R$ and for all ideals J with $I \subseteq J \subseteq R$, $J = I$ or $J = R$. All maximal ideals are prime, but not all prime ideals are maximal.
3. *Prime and Irreducible Elements.* Given x in a commutative ring R , x is *prime in R* iff $x \neq 0$, x is not a unit of R , and for all $y, z \in R$, $x|(yz)$ implies $x|y$ or $x|z$; x is *irreducible in R* iff $x \neq 0$, x is not a unit of R , and for all $w \in R$, if $w|x$ then $w \sim x$ or w is a unit. In terms of ideals, x is prime iff Rx is a nonzero prime ideal of R ; x is irreducible iff Rx is a nonzero ideal that is maximal in the poset of proper, principal ideals of R . In an integral domain, every prime element is irreducible; the converse fails in general but holds in PIDs and UFDs.
4. *Types of Rings.* Let R be a commutative ring. R is an *integral domain* iff $0_R \neq 1_R$ and R has no (nonzero) zero divisors. R is a *principal ideal domain (PID)* iff R is an integral domain in which every ideal I has the form $I = Rc$ for some $c \in R$. R is a *unique factorization domain (UFD)* iff R is an integral domain in which (a) every $r \neq 0$ in R factors as $r = up_1 \cdots p_k$ for some $u \in R^*$ and irreducible $p_i \in R$; and (b) whenever $up_1 \cdots p_k = vq_1 \cdots q_t$ with $u, v \in R^*$ and all p_i, q_j irreducible in R , $k = t$ and $p_i \sim q_i$ for $1 \leq i \leq k$ after reordering the q_j 's.
5. *Theorems about PIDs and UFDs.* \mathbb{Z} and $F[x]$ (for any field F) are PIDs. For all x in a PID or UFD R , x is prime in R iff x is irreducible in R . Every PID is a UFD. Given a_1, \dots, a_k in a PID R , there exist gcds and lcms of a_1, \dots, a_k , and each gcd is an R -linear combination of the a_j 's. More precisely, d is a gcd of a_1, \dots, a_k iff $Rd = Ra_1 + \cdots + Ra_k$, and e is an lcm of a_1, \dots, a_k iff $Re = Ra_1 \cap \cdots \cap Ra_k$. Every ascending chain of ideals in a PID must stabilize. Every nonzero prime ideal in a PID is maximal. In a UFD R , gcds and lcms of a_1, \dots, a_k exist, but the gcd may not be an R -linear combination of the a_i 's (if it is in all cases, then R must be a PID).
6. *Module Lemmas.* Assume R is a ring and M and N are isomorphic R -modules. If R is *commutative*: for any $c \in R$, cM is a submodule of M with $cM \cong cN$ and $M/cM \cong N/cN$; every basis of an f.g. free R -module has the same size; and the annihilator $\text{ann}_R(M) = \{r \in R : \forall x \in M, rx = 0\}$ is an ideal and equals $\text{ann}_R(N)$. If R is an *integral domain*: the torsion module $\text{tor}(M) = \{x \in M : \exists r \in R \sim \{0\}, rx = 0\}$ is a submodule of M with $\text{tor}(M) \cong \text{tor}(N)$ and $M/\text{tor}(M) \cong N/\text{tor}(N)$. If R is a *PID*: every submodule of an f.g. free R -module is also f.g. free; given $a = p_1^{e_1} \cdots p_k^{e_k} \in R$ where the p_i are non-associate irreducible elements of R , $R/Ra \cong \prod_{i=1}^k R/Rp_i^{e_i}$.

7. *Smith Normal Form of a Matrix.* For a PID R and any matrix $A \in M_{m,n}(R)$, there exist unique $s \in \mathbb{N}$ and unique ideals $Ra_1 \supseteq Ra_2 \supseteq \dots \supseteq Ra_s \neq \{0\}$ such that for some invertible matrices $P \in M_m(R)$ and $Q \in M_n(R)$, $PAQ = \text{diag}(a_1, \dots, a_s)_{m \times n}$. Note $a_i | a_{i+1}$ for $1 \leq i < s$. We say s is the *rank* of A , the ideals Ra_i are the *invariant factors* of A , and the matrix PAQ is a *Smith normal form* of A . We can pass from A to PAQ using a sequence of type 4 row and column operations. For $1 \leq k \leq s$, $a_1 a_2 \dots a_k$ is a gcd of all the k 'th order minors of A (which are determinants of $k \times k$ submatrices obtained by keeping any k rows and any k columns of A), and s is the largest k for which some k 'th order minor of A is nonzero.
8. *Structure Theorem for Linear Maps between Free Modules.* Let N and M be f.g. free R -modules where R is a PID. For every R -linear map $T : N \rightarrow M$, there exist unique $s \in \mathbb{N}$ and unique ideals $Ra_1 \supseteq Ra_2 \supseteq \dots \supseteq Ra_s \neq \{0\}$ such that for some ordered basis $X = (x_1, \dots, x_n)$ of N and some ordered basis $Y = (y_1, \dots, y_m)$ of M , $T(x_i) = a_i y_i$ for $1 \leq i \leq s$ and $T(x_i) = 0_M$ for $s < i \leq n$.
9. *Structure of Finitely Generated Modules over PIDs.* For any finitely generated module M over a PID R , there exist unique $k, d \in \mathbb{N}$ and a unique sequence of ideals $R \neq Ra_1 \supseteq Ra_2 \supseteq \dots \supseteq Ra_k \neq \{0\}$ such that $M \cong R/Ra_1 \times R/Ra_2 \times \dots \times R/Ra_k \times R^d$. There also exist unique $m, d \in \mathbb{N}$ and a unique multiset of ideals $[Rq_1, \dots, Rq_m]$, with every q_i a power of an irreducible in R , such that $M \cong R/Rq_1 \times \dots \times R/Rq_m \times R^d$. We call d the *Betti number* of M , (Ra_1, \dots, Ra_k) the *invariant factors* of M , and $[Rq_1, \dots, Rq_m]$ the *elementary divisors* of M .
10. *Rational Canonical Forms.* Let F be any field. Given $h = x^d + h_{d-1}x^{d-1} + \dots + h_0 \in F[x]$, the *companion matrix* C_h has $C_h(i+1, i) = 1$ for $1 \leq i < d$, $C_h(i, d) = -h_{i-1}$ for $1 \leq i \leq d$, and all other entries zero. A *rational canonical form* is a matrix of the form blk-diag(C_{h_1}, \dots, C_{h_k}), where $h_1, \dots, h_k \in F[x]$ are monic and non-constant and $h_i | h_{i+1}$ for $1 \leq i < k$. For every n -dimensional vector space V and every linear map $T : V \rightarrow V$, there exists a unique rational canonical form $B \in M_n(F)$ such that $[T]_X = B$ for some ordered basis X of V . For any matrix $A \in M_n(F)$, there exists a unique rational canonical form $B \in M_n(F)$ similar to A (i.e., $B = P^{-1}AP$ for some invertible $P \in M_n(F)$). In this case, T and A and B have minimal polynomial h_k and characteristic polynomial $h_1 \dots h_k$, so $m_A | \chi_A$. The invariant factors h_1, \dots, h_k in the rational canonical form of $A \in M_n(F)$ coincide with the non-constant monic invariant factors in the Smith normal form of $xI_n - A \in M_n(F[x])$.
11. *Jordan Canonical Forms.* Let F be an algebraically closed field (such as \mathbb{C}). For $c \in F$, the *Jordan block* $J(c; e)$ is the $e \times e$ matrix with c 's on the main diagonal, 1's on the next higher diagonal, and zeroes elsewhere. A *Jordan canonical form* is a matrix of the form blk-diag($J(c_1; e_1), \dots, J(c_s; e_s)$). For every n -dimensional F -vector space V and every linear map $T : V \rightarrow V$, there exists a Jordan canonical form $C \in M_n(F)$ such that $[T]_X = C$ for some ordered basis X of V . For any matrix $A \in M_n(F)$, there exists a Jordan canonical form $C \in M_n(F)$ similar to A . The only other Jordan forms that can appear here are obtained by reordering the Jordan blocks of C .
12. *$F[x]$ -Module of a Linear Operator.* Given a field F , a finite-dimensional F -vector space V , and a linear map $T : V \rightarrow V$, V becomes an $F[x]$ -module via $(\sum_{i=0}^d p_i x^i) \cdot v = \sum_{i=0}^d p_i T^i(v)$ for $p_i \in F$ and $v \in V$. $F[x]$ -submodules of V are the *T -invariant subspaces* (subspaces W with $T[W] \subseteq W$). A cyclic $F[x]$ -

submodule W has an F -basis of the form $B_z = (z, T(z), T^2(z), \dots, T^{d-1}(z))$ for some $z \in V$ and $d \in \mathbb{N}^+$. Such a cyclic submodule is isomorphic to $F[x]/F[x]h$, where $h \in F[x]$ is the monic polynomial of least degree d with $h \cdot z = 0$. The matrix of $[T|_W]_{B_z}$ is the companion matrix C_h .

18.21 Exercises

1. (a) Prove: for all commutative rings R and all $a \in R$, Ra is an ideal of R . (b) Give an example to show (a) can be false if R is not commutative.
2. Prove that every field F is a PID.
3. (a) Let $R = \mathbb{Z}[x]$. Prove that $I = \{2f + xg : f, g \in R\}$ is an ideal of R that is not principal. Conclude that R is not a PID. (b) For any integral domain R and integer $n \geq 2$, prove $S = R[x_1, \dots, x_n]$ is not a PID.
4. Let $R = \{a + bi : a, b \in \mathbb{Z}\}$. (a) Prove R is a subring of \mathbb{C} and hence an integral domain. (b) Prove: for all $f, g \in R$ with $g \neq 0$, there exist $q, r \in R$ with $f = qg + r$ and $r = 0$ or $|r| < |g|$, where $|a + bi| = \sqrt{a^2 + b^2}$. [Hint: First modify the division algorithm in \mathbb{Z} to see that for all $u, v \in \mathbb{Z}$ with $v \neq 0$, there exist $q, r \in \mathbb{Z}$ with $u = qv + r$ and $|r| \leq v/2$.] (c) Use (a) and (b) to prove that R is a PID.
5. A *Euclidean domain* is a ring R satisfying the following hypotheses. First, R is an integral domain. Second, R has a *degree function* $\deg : R \setminus \{0_R\} \rightarrow \mathbb{N}$ such that for all nonzero $f, g \in R$, $\deg(f) \leq \deg(fg)$. Third, R satisfies a *division theorem* relative to the degree function: for all $f, g \in R$ with $g \neq 0$, there exist $q, r \in R$ with $f = qg + r$ and $r = 0$ or $\deg(r) < \deg(g)$. (a) Prove that every Euclidean domain is a PID. (b) Explain how to use (a) to show that fields F , polynomial rings $F[x]$, and the ring \mathbb{Z} are PIDs. [It can be shown [13] that not all PIDs are Euclidean domains; the standard example is the ring $R = \{a + bz : a, b \in \mathbb{Z}\}$, where $z = (1 + i\sqrt{19})/2$.]
6. Let F be a field. (a) Prove: $F[[x]]^*$ consists of all formal power series $\sum_{i \geq 0} p_i x^i$ with $p_0 \neq 0$. (b) Prove: for any field F , $F[[x]]$ is a PID. [Hint: Show that every nonzero ideal of $F[[x]]$ has the form $F[[x]]x^j$ for some $j \in \mathbb{N}$.]
7. Let R be a commutative ring and fix $a, a', b, c \in R$. Prove: (a) $a|a$; (b) if $a|b$ and $b|c$ then $a|c$; (c) $1_R|a$ and $a|0_R$; (d) $0|a$ iff $a = 0$; (e) \sim (defined by $a \sim b$ iff $a|b$ and $b|a$) is an equivalence relation on R ; (f) if $a \sim a'$, then $a|b$ iff $a'|b$ and $c|a$ iff $c|a'$; (g) $a \in R^*$ iff $a \sim 1$.
8. Let R be a commutative ring. (a) Prove: for all $k \in \mathbb{N}^+$, and all $r_1, \dots, r_k, a_1, \dots, a_k, c \in R$, if c divides every r_i , then $c|(a_1r_1 + \dots + a_kr_k)$. (b) Translate the result in (a) into a statement about principal ideals.
9. Suppose R is a commutative ring and $d \in R$ is a gcd of $a_1, \dots, a_k \in R$. (a) Prove the set of all gcds of a_1, \dots, a_k in R equals the set of all $d' \in R$ with $d \sim d'$. (b) Prove a similar result for lcms.
10. Let $R = \{a + bi : a, b \in \mathbb{Z}\}$, which is a subring of \mathbb{C} . (a) Find all units of R . (b) What are the associates of $3 - 4i$ in R ? (c) How do the answers to (a) and (b) change if we replace R by \mathbb{C} ?
11. Given d, a_1, \dots, a_k in a commutative ring R , carefully prove that d is an lcm of

a_1, \dots, a_k iff Rd is the greatest lower bound of $\{Ra_1, \dots, Ra_k\}$ in the poset of principal ideals of R ordered by \subseteq .

12. Suppose R is a commutative ring and $a, b \in R$ are associates. (a) Prove a is irreducible in R iff b is irreducible in R using the definitions involving divisibility of ring elements. (b) Prove the result in (a) by considering principal ideals. (c) Prove a is prime in R iff b is prime in R using divisibility definitions. (d) Prove the result in (c) by considering principal ideals.
13. (a) Suppose p is a prime element in a commutative ring R . Prove by induction: for all $k \in \mathbb{N}^+$ and all $a_1, \dots, a_k \in R$, if $p|(a_1a_2 \cdots a_k)$ then p divides some a_i . (b) Suppose p is irreducible in an integral domain R . Prove: for all $k \in \mathbb{N}^+$ and all $b_1, \dots, b_k \in R$, if $p \sim b_1b_2 \cdots b_k$, then for some i , $p \sim b_i$ and all other b_j 's are units of R .
14. Let R be a commutative ring. (a) Prove an ideal P of R is a prime ideal iff R/P is an integral domain. (b) Prove an ideal M of R is a maximal ideal iff R/M is a field. (c) Deduce from (a) and (b) that every maximal ideal is a prime ideal.
15. Let R be a commutative ring with a maximal ideal M . Without using quotient rings, prove that M is a prime ideal. [Hint: Arguing by contradiction, assume $x, y \in R$ satisfy $x, y \notin M$ but $xy \in M$. Consider the ideals $M + Rx$ and $M + Ry$ to deduce $1_R \in M$, which cannot occur.]
16. Suppose R is a PID and $I = Ra$ is an ideal of R . Under what conditions on a is R/I a PID? Explain.
17. Consider the commutative ring $R = \mathbb{Z}_{12}$, which is not an integral domain. (a) Find all units of \mathbb{Z}_{12} . (b) Find all prime elements in \mathbb{Z}_{12} . (c) Find all irreducible elements in \mathbb{Z}_{12} . (d) Describe the equivalence classes of \sim in \mathbb{Z}_{12} . (e) Is it true that for all $a, b \in \mathbb{Z}_{12}$, $a \sim b$ iff $a = ub$ for some $u \in \mathbb{Z}_{12}^*$?
18. Let R and S be rings. Suppose K is any ideal in the product ring $R \times S$. Prove there exists an ideal I in R and an ideal J in S such that $K = I \times J$.
19. Consider the product ring $T = \mathbb{Z}_4 \times \mathbb{Z}_2$. (a) Find all ideals in T . Draw the ideal lattice of T (compare to Figure 17.1). Decide, with explanation, whether each ideal is prime and whether each ideal is maximal. (b) Aided by (a), decide (with explanation) whether each nonzero element of T is a unit, a zero divisor, a prime element, or an irreducible element (indicate all that apply).
20. Repeat the previous exercise for the product ring $T = \mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_2$. (Compare the ideal lattice to the one in Exercise 45 of Chapter 17.)
21. Let R be a PID that is not a field. (a) Explain why $\{0\}$ is a prime ideal of R that is not a maximal ideal. (b) Show that every nonzero prime ideal of R is a maximal ideal.
22. Suppose R is an integral domain in which every ascending chain of *principal* ideals stabilizes. Prove that the existence assertion (a) in the definition of a UFD (§18.5) holds in R .
23. Suppose R is an integral domain in which every irreducible element is prime. Prove that the uniqueness assertion (b) in the definition of a UFD (§18.5) holds in R .
24. Suppose R is a UFD. Prove $p \in R$ is irreducible in R iff p is prime in R .
25. Let R be a UFD, and let $\{p_i : i \in I\}$ be a set of irreducible elements of R such that every irreducible q in R is associate to exactly one p_i . (a) Prove: for all

nonzero $r \in R$, there exist unique $u \in \mathbb{R}^*$ and exponents $e_i \in \mathbb{N}$ such that all but finitely many e_i are zero and $r = u \prod_{i \in I} p_i^{e_i}$. Write $u = u(r)$ and $e_i = e_i(r)$ to indicate that these parameters are functions of r . (b) Prove: for all $a, b, c \in R$, $c = ab$ iff $u(c) = u(a)u(b)$ and $e_i(c) = e_i(a) + e_i(b)$ for all $i \in I$. (c) Prove: for all nonzero $r, s \in R$, $r|s$ iff $e_i(r) \leq e_i(s)$ for all $i \in I$. (d) Prove: for all nonzero $d, a_1, \dots, a_k \in R$, d is a gcd of a_1, \dots, a_k iff $e_i(d) = \min_{1 \leq j \leq k} e_i(a_j)$ for all $i \in I$. Deduce that *gcds of finite lists of elements in a UFD always exist*. (e) State and prove a similar formula characterizing lcms of a_1, \dots, a_k .

26. Suppose R is a Euclidean domain (Exercise 5) for which there is an algorithm to compute the quotient and remainder (relative to deg) when f is divided by g . Describe an algorithm that takes as input $f, g \in R$ and produces as output $d, a, b \in R$ such that $d = af + bg$ and d is a gcd of f and g . (Imitate the proof in §3.8.)
27. Let $R = \{a + bi\sqrt{5} : a, b \in \mathbb{Z}\}$, which is a subring of \mathbb{C} . (a) Define $N : R \rightarrow \mathbb{Z}$ by $N(a + bi\sqrt{5}) = a^2 + 5b^2$ for $a, b \in \mathbb{Z}$. Prove $N(rs) = N(r)N(s)$ for all $r, s \in R$. (b) Prove: for all $r \in R$, r is a unit of R iff $N(r) = \pm 1$, and find all units of R . (c) Show that $2, 3, 1 + i\sqrt{5}$, and $1 - i\sqrt{5}$ are irreducible in R but not prime in R . (d) Is R a UFD? Is R a PID? Why?
28. Let $R = \mathbb{Z}_6$. Show that every ideal of R is principal, but give an example of a f.g. free R -module M and a submodule N of M that is not free. Explain exactly why the proof in §18.6 fails for this ring.
29. Give a specific example of an integral domain R , a f.g. free R -module M , and a submodule N of M that is not free. Why does the proof in §18.6 fail here?
30. Prove or disprove: for every f.g. free module M over a PID R and every submodule N of M , there is a submodule P of M with $P + N = M$ and $P \cap N = \{0\}$.
31. Without using operation (B4), give direct proofs that applying the operations (B1), (B2), and (B3) in §18.7 to an ordered basis of M produces a new ordered basis of M .
32. In §18.8: (a) prove claim 1; (b) use claims 1 and 2 to deduce statements 3, 4, and 5; (c) prove properties 6, 7, and 8.
33. Let $R = \mathbb{Q}[x]$, let $X = (e_1, e_2) = Y$ be the standard ordered basis of the R -module R^2 , and let $T : R^2 \rightarrow R^2$ be the R -linear map with matrix $A = \begin{bmatrix} x-3 & x^2+1 \\ x^3-x & 5x+2 \end{bmatrix}$ relative to the bases X and Y . (a) Compute $T((3x-1, x^2+x-1))$. (b) Let $Z = ((x-1, 2x), (x/2, x+1))$. Show Z is obtained from X by operation (B4), so is an ordered basis of R^2 . (c) What is the matrix of T relative to input basis Z and output basis Y ? How is this matrix related to A ? (d) What is the matrix of T relative to input basis X and output basis Z ? How is this matrix related to A ? (e) What is the matrix of T relative to input basis Z and output basis Z ? How is this matrix related to A ? (f) Find a single type 4 column operation on A that will make the 1, 1-entry become 1. Find X' such that the new matrix represents T relative to X' and Y . (g) Find a single type 4 row operation on A that will make the 1, 1-entry become 1. Find Y' such that the new matrix represents T relative to X and Y' .
34. Let M be a f.g. free module over a PID R , let $X = (x_1, \dots, x_m)$ be an ordered basis of M , and let $P \in M_m(R)$ be an invertible matrix. Consider the operation that maps X to $X' = (x'_1, \dots, x'_m)$, where $x'_i = \sum_{j=1}^m P(j, i)x_j$ for $1 \leq i \leq m$. (a) Prove X' is an ordered basis of M . (b) Prove that for any ordered basis Y of

M , we can choose $P \in M_m(R)$ such that $X' = Y$. (c) Show that operation (B4) is a special case of the operation considered here.

35. Let R be a PID, and let $T : N \rightarrow M$ be an R -linear map between two f.g. free R -modules. Let N and M have respective ordered bases $X = (x_1, \dots, x_n)$ and $Y = (y_1, \dots, y_m)$, and let A be the matrix of T relative to X and Y . (a) Given an invertible $P \in M_n(R)$, define $X' = XP$ as in Exercise 34, and let A' be the matrix of T relative to X' and Y . How is A' related to A and P ? (b) Given an invertible $Q \in M_m(R)$, let $Y' = YQ$ as in Exercise 34, and let A'' be the matrix of T relative to X and Y' . How is A'' related to A and Q ? (c) Suppose X' and Y' are any ordered bases for N and M (respectively), and B is the matrix of T relative to X' and Y' . Show that $B = UAV$ for some invertible $U \in M_m(R)$ and $V \in M_n(R)$.
36. Let a, b, d be nonzero elements in a PID (or UFD) R . Prove the following assertions from §18.9: (a) $\text{len}(ab) = \text{len}(a) + \text{len}(b)$; (b) if $d|a$ in R , then $\text{len}(d) \leq \text{len}(a)$ with equality iff $d \sim a$; (c) if $d = \gcd(a, b)$ and a does not divide b , then $\text{len}(d) < \text{len}(a)$.
37. Suppose R is a UFD such that for all $x, y \in R$ and every gcd d of x and y in R , there exist $a, b \in R$ with $d = ax + by$. Prove that R must be a PID. (This means that in UFDs that are not PIDs, the gcd of a list of elements is *not* always an R -linear combination of those elements.) [Hint for the proof: given any nonzero ideal I of R , show I is generated by any $x \in I$ of minimum length.]
38. Prove the matrix reduction theorem in §18.9 for Euclidean domains R (see Exercise 5) without using unique factorization or length, by imitating the proof for $R = \mathbb{Z}$ in §15.11. Also show that the matrix reduction never requires general type 4 operations, but can be achieved using only elementary operations (B1), (B2), and (B3).
39. Let R be the PID $\{a + bi : a, b \in \mathbb{Z}\}$ (see Exercise 4). Use matrix reduction to find a Smith normal form of these matrices: (a) $\begin{bmatrix} 8 - 21i & 12 + 51i \\ 49 + 32i & -136 + 2i \end{bmatrix}$; (b) $\begin{bmatrix} 4+i & 3 & 1+i \\ 16+5i & 10-7i & 1-3i \\ 1+10i & 8-3i & 3-7i \end{bmatrix}$.
40. Solve Exercise 39 again, but find the rank and invariant factors by computing gcds of k 'th order minors.
41. Let $R = \mathbb{Q}[x]$. Use matrix reduction to find a Smith normal form of these matrices: (a) $\begin{bmatrix} x^3 - x^2 - 3x + 2 & x^2 - 4 \\ x^3 - 6x + 4 & x^2 + x - 6 \end{bmatrix}$; (b) $\begin{bmatrix} x^4 - x^2 + x & x^4 + x^3 - x^2 & x^3 - x \\ x^5 - x^4 - x^3 + 2x^2 - 2x & x^5 - 2x^3 + x^2 - x & x^4 - x^3 - x^2 + x \end{bmatrix}$; (c) $\begin{bmatrix} x+1 & x^2 + x - 1 \\ 1 & x \\ x^2 + x & x^3 + x^2 - x \end{bmatrix}$.
42. Solve Exercise 41 again, but find the rank and invariant factors by computing gcds of k 'th order minors.
43. (a) Let R be a PID and $a \in R$. What is a Smith normal form of an $m \times n$ matrix all of whose entries equal a ? (b) Given $c \in \mathbb{C}$, what is a Jordan canonical form of an $n \times n$ complex matrix all of whose entries equal c ? (c) Given $c \in \mathbb{Q}$, what is the rational canonical form of the matrix in $M_n(\mathbb{Q})$ all of whose entries equal c ?

- (d) For prime p , what is the rational canonical form of the matrix in $M_p(\mathbb{Z}_p)$ all of whose entries equal 1?
44. (a) Write a computer program that finds the Smith normal form of a matrix $A \in M_{m,n}(\mathbb{Q}[x])$ via the matrix reduction algorithm described in §18.9. (b) Write another program to find the Smith normal form of A using the formulas for the invariant factors as gcds of minors of A (see §18.12). (c) Comment on the relative efficiency of the programs in (a) and (b) for large m and n .
45. Let $R = \mathbb{Q}[x]$, and let $T : R^3 \rightarrow R^2$ be the R -linear map
- $$T((f, g, h)) = (x^3 f + (x^2 + x)g + x^2 h, (x^4 - x^3)f + x^2 g + x^3 h) \text{ for } f, g, h \in R.$$
- (a) What is the matrix of T relative to the standard ordered bases of R^3 and R^2 ?
 (b) Find a Smith normal form B of T and bases X of R^3 and Y of R^2 such that B is the matrix of T relative to X and Y .
46. For each R -module M , find an isomorphic module of the form (18.5) and of the form (18.6). (a) $R = \mathbb{Q}[x]$, $M = R^4/(Rv_1 + Rv_2 + Rv_3)$ where
- $$\begin{aligned} v_1 &= (x^4 - 2x^3 + 2x^2, -2x^2, -x^2, x^2), \\ v_2 &= (x^3 - x^2, x^2, x^2, 0), \\ v_3 &= (x^5 - 2x^3 + 2x^2, x^4 - 2x^2, x^4 - x^2, x^2). \end{aligned}$$
- (b) $R = \mathbb{Q}[x]$, $M = R^3/(Rw_1 + Rw_2 + Rw_3)$ where $w_1 = (x^2, x^3, x^3)$, $w_2 = (x^3, x^3, x^2)$, and $w_3 = (x^2, x^3, x^4)$. (c) $R = \{a+bi : a, b \in \mathbb{Z}\}$, $M = R^2/(Rz_1 + Rz_2)$ where $z_1 = (4 + 3i, 12 + 8i)$ and $z_2 = (5 - 7i, 7 - 4i)$.
47. (a) Let R be a PID. Prove: if $a, b \in R$ satisfy $\gcd(a, b) = 1$, then R/Rab is isomorphic to $R/Ra \times R/Rb$ both as R -modules and as rings. (b) Prove that (a) need not hold when R is a UFD that is not a PID by showing that for $R = \mathbb{Q}[x, y]$, R/Rxy and $R/Rx \times R/Ry$ are not isomorphic as R -modules.
48. Let A be an $m \times n$ matrix with entries in a PID R . (a) Use the Cauchy–Binet formula to prove: for all $1 \leq k \leq \min(m, n)$ and all invertible $Q \in M_n(R)$, $g_k(AQ) \sim g_k(A)$. (b) If Q is not invertible, is there any relation between $g_k(AQ)$ and $g_k(A)$?
49. Let A be an $m \times n$ matrix with entries in a PID R . Without using the Cauchy–Binet formula, prove: for all $1 \leq k \leq \min(m, n)$ and all invertible $P \in M_m(R)$ and all invertible $Q \in M_n(R)$, $g_k(PA) \sim g_k(A) \sim g_k(AQ)$. [Use ideas from §4.8 to show that each row of $(PA)_{I,J}$ is an R -linear combination of certain rows of $A_{[m],J}$. Then use multilinearity of the determinant as a function of the rows of a matrix (§5.6) to show that $\det((PA)_{I,J})$ is some R -linear combination of k 'th order minors of A . The Cauchy–Binet formula shows us explicitly what this linear combination is.]
50. Give an example of a commutative ring R and an R -module M such that $\text{tor}(M)$ is not a submodule of M .
51. Suppose R is an integral domain and $f : M \rightarrow N$ is an R -module isomorphism. Carefully check that $f[\text{tor}(M)] = \text{tor}(N)$ and that f induces an isomorphism $f' : M/\text{tor}(M) \rightarrow N/\text{tor}(N)$.
52. Let a be a nonzero element in a PID R . Prove that $\text{len}(R/Ra)$ [as defined in §17.11] equals $\text{len}(a)$ [as defined in §18.9].

53. Let N and N' be isomorphic R -modules, where R is a commutative ring, and fix $c \in R$. (a) Prove cN is a submodule of N . (b) Prove $cN \cong cN'$ and $N/cN \cong N'/cN'$. (c) If $N = N_1 \times \cdots \times N_k$, prove $cN = cN_1 \times \cdots \times cN_k$.
54. Define $T : \mathbb{R}^4 \rightarrow \mathbb{R}^4$ by $T(v) = Av$ for $v \in \mathbb{R}^4$, where $A = \text{blk-diag}(J(3; 2), J(5; 2))$ (a Jordan canonical form). (a) What is the rational canonical form of T ? (b) Find all T -invariant subspaces of \mathbb{R}^4 . Find a specific generator for each T -cyclic subspace. [Hint: Use (a) and results about submodule lattices.]
55. (a) Repeat Exercise 54(a), taking $A = \text{blk-diag}(J(4; 2), J(4; 2))$. (b) For $T(x) = Ax$, show that \mathbb{R}^4 has infinitely many 1-dimensional T -invariant subspaces and infinitely many 2-dimensional T -invariant subspaces.
56. (a) A certain matrix in $M_n(\mathbb{Q})$ has invariant factors $(x - 1, x^3 - 3x + 2, x^5 + x^4 - 5x^3 - x^2 + 8x - 4)$. Find n and the elementary divisors of the matrix. (b) A certain matrix in $M_m(\mathbb{Q})$ has elementary divisors $[x - 1, x - 1, x - 1, (x - 1)^3, (x - 1)^4, (x^2 - 2)^2, (x^2 - 2)^2, (x^2 - 2)^3, x^2 + 1, x^2 + 1]$. Find m and the invariant factors of the matrix.
57. Let M be a finitely generated module over a PID R . Give an alternate proof of the uniqueness of the elementary divisors of M by imitating the arguments for \mathbb{Z} -modules in §15.17 and §15.18.
58. (a) In §18.15, prove $f(g(L)) = L$ for all $L \in Y$. (b) Let R be a PID. Assume we have proved the uniqueness of the elementary divisors of a finitely generated R -module M (see Exercise 57). Use this result and (a) to prove the uniqueness of the invariant factors of M (cf. §15.19).
59. In §18.16, we proved that V was an $F[x]$ -module by checking all the module axioms. Give a more conceptual proof of this fact by using Exercise 12 in Chapter 17 and the universal mapping property for $F[x]$.
60. Let $T, S : V \rightarrow V$ be linear operators on an n -dimensional vector space V over a field F with $S \circ T = T \circ S$. Show that there exists a unique $F[x, y]$ -module action $\cdot : F[x, y] \times V \rightarrow V$ such that for all $v \in V$, $x \cdot v = T(v)$, $y \cdot v = S(v)$, and for all $c \in F$, $c \cdot v$ is the given scalar multiplication in V .
61. Let F be a field, and let $h \in F[x]$ have degree $d > 0$. Show that the F -vector space $F[x]/F[x]h$ has ordered basis $(1 + F[x]h, x + F[x]h, \dots, x^{d-1} + F[x]h)$.
62. Let $T : \mathbb{R}^6 \rightarrow \mathbb{R}^6$ be given by $T((x_1, x_2, x_3, x_4, x_5, x_6)) = (x_2, x_3, x_4, x_5, x_6, x_1)$ for $x_i \in \mathbb{R}$. (a) For $1 \leq i \leq 6$, find the T -annihilator of $e_i \in \mathbb{R}^6$ in $\mathbb{R}[x]$. (b) What is the T -annihilator of $e_1 + e_3 + e_5$? (c) What is the T -annihilator of $e_1 + e_4$? (d) If possible, find $v \in \mathbb{R}^6$ whose T -annihilator is $\mathbb{R}[x](x^2 - x + 1)$. (e) If possible, find $w \in \mathbb{R}^6$ whose T -annihilator is $\mathbb{R}[x](x^3 + 2x^2 + 2x + 1)$.
63. (a) Prove the assertion in the last sentence of §18.16. (b) In the last paragraph of §18.17, confirm that each W_i^* is a cyclic $F[x]$ -submodule of V generated by z_i^* , and $V \cong W_1^* \times \cdots \times W_\ell^* \cong \prod_{i=1}^\ell F[x]/F[x]g_i$.
64. In §18.18, check carefully that $B = (y_{e-1}, \dots, y_2, y_1, y_0)$ is an ordered F -basis for W .
65. (a) In §18.18, prove that if W is a T -invariant subspace of V such that $[T|_W]_B = J(c; e)$ for some ordered basis B of W , then W is a T -cyclic subspace with $W = F[x]y \cong F[x]/F[x](x - c)^e$ for some $y \in W$. (b) Use (a) to show that all Jordan canonical forms of T are obtained from (18.14) by reordering the Jordan blocks.

66. Let $h \in F[x]$ be monic of degree $n > 0$. (a) Show the minimal polynomial of C_h is h . (b) Show that $xI_n - C_h \in M_n(F[x])$ can be reduced by elementary row and column operations to $\text{diag}(1, \dots, 1, h)$. [Hint: Start by adding $-x$ times row i to row $i-1$, for $i = n, n-1, \dots, 2$.] (c) Show that $\chi_{C_h} = h$ by calculating $\det(xI_n - C_h)$. [Hint: You can use induction on n , or analyze your answer to (b).]
67. (a) Show that if $A \in M_n(F)$ is similar to B , then $m_A = m_B$ and $\chi_A = \chi_B$. (b) Show that for a block-diagonal matrix $B = \text{blk-diag}(B_1, \dots, B_k)$, $m_B = \text{lcm}(m_{B_1}, \dots, m_{B_k})$ and $\chi_B = \prod_{i=1}^k \chi_{B_i}$.
68. (a) What are the minimal polynomial and characteristic polynomial of a Jordan block $J(c; e)$? (b) What are the minimal polynomial and characteristic polynomial of a Jordan canonical form $\text{blk-diag}(J(c_1; e_1), \dots, J(c_s; e_s))$?
69. Let F be a field. (a) What are the invariant factors in the rational canonical form of the matrix bI_n ? (b) What are the invariant factors for a diagonal matrix $\text{diag}(a_1, \dots, a_n)_{n \times n}$ where all $a_i \in F$ are distinct? (c) Given that $a, b, c, d \in F$ are distinct, what are the invariant factors of $\text{diag}(a, a, a, a, b, b, b, b, c, c, d, d)$? (d) Find (with proof) all matrices in $M_n(F)$ whose rational canonical form is diagonal.
70. (a) Use the Jordan canonical form to prove that a matrix $A \in M_n(\mathbb{C})$ is diagonalizable iff the minimal polynomial of A in $\mathbb{C}[x]$ has no repeated roots. (b) For any field F , give a simple characterization of the rational canonical forms of diagonalizable matrices in $M_n(F)$.
71. Find the rational canonical form of the following matrices in $M_n(\mathbb{Q})$:
- (a) $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$; (b) $\begin{bmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$; (c) $\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 & -1 & 0 \\ 0 & 4 & 3 & 0 \\ 1 & -4 & -2 & 1 \end{bmatrix}$;
- (d) $\begin{bmatrix} 12 & 15 & 5 & -30 \\ 1 & 7/2 & 1/2 & -3 \\ 13 & 39/2 & 17/2 & -39 \\ 6 & 9 & 3 & -16 \end{bmatrix}$
72. Find a Jordan canonical form in $M_n(\mathbb{C})$ for each matrix in Exercise 71.
73. (a) Prove: for any PID R and any $B \in M_{m,n}(R)$, B and B^T have the same rank and invariant factors. (b) Prove: for any field F and all $A \in M_n(F)$, A is similar to A^T .
74. Let F be a subfield of K . Prove: for all $A \in M_n(F)$, the rational canonical form of A in $M_n(F)$ is the same as the rational canonical form of A in $M_n(K)$.

Part VI

Universal Mapping Properties and Multilinear Algebra

This page intentionally left blank

Introduction to Universal Mapping Properties

The concept of a *Universal Mapping Property* (abbreviated UMP) occurs frequently in linear and multilinear algebra. Indeed, this concept pervades every branch of abstract algebra and occurs in many other parts of mathematics as well. To introduce this fundamental notion, we start by considering a well-known result about vector spaces over a field F .

Theorem: Let V be an n -dimensional F -vector space with basis $X = \{x_1, \dots, x_n\}$. For any F -vector space W and any function $f : X \rightarrow W$, there exists a unique F -linear map $T : V \rightarrow W$ extending f (i.e., $T(x) = f(x)$ for all $x \in X$).

Proof of uniqueness of T : Any $v \in V$ can be uniquely expressed as $v = \sum_{i=1}^n c_i x_i$ for scalars $c_i \in F$. If T is to be an F -linear map extending f , then we are forced to have

$$T(v) = T\left(\sum_{i=1}^n c_i x_i\right) = \sum_{i=1}^n c_i T(x_i) = \sum_{i=1}^n c_i f(x_i). \quad (19.1)$$

This shows that the map T is uniquely determined by f , if T exists at all.

Proof of existence of T : To prove existence of T , we define $T(\sum_{i=1}^n c_i x_i) = \sum_{i=1}^n c_i f(x_i)$ for all $c_i \in F$, as in formula (19.1). Note that T is a well-defined map from V into W , since every vector $v \in V$ has a unique expansion in terms of the basis X . Taking $v = x_j = 1x_j + \sum_{i \neq j} 0x_i$, the formula shows that $T(x_j) = 1f(x_j) + \sum_{i \neq j} 0f(x_i) = f(x_j)$ for all $x_j \in X$, so that T does extend f . To check F -linearity of T , fix $v, w \in V$ and $a \in F$. Write $v = \sum_{i=1}^n c_i x_i$ and $w = \sum_{i=1}^n d_i x_i$ for some $c_i, d_i \in F$. Then $v + w = \sum_{i=1}^n (c_i + d_i)x_i$ and $av = \sum_{i=1}^n (ac_i)x_i$. Using the definition of T several times, we get

$$\begin{aligned} T(v + w) &= \sum_{i=1}^n (c_i + d_i)f(x_i) = \sum_{i=1}^n c_i f(x_i) + \sum_{i=1}^n d_i f(x_i) = T(v) + T(w); \\ T(av) &= \sum_{i=1}^n (ac_i)f(x_i) = a \sum_{i=1}^n c_i f(x_i) = aT(v). \end{aligned}$$

Therefore, T is an F -linear map.

We now restate the theorem just proved in three equivalent ways. Throughout, we fix the field F , the vector space V , and the basis X . Let $i : X \rightarrow V$ be the inclusion map given by $i(x) = x \in V$ for all $x \in X$. Note that T extends f iff $T(x) = f(x)$ for all $x \in X$ iff $T(i(x)) = f(x)$ for all $x \in X$ iff $(T \circ i)(x) = f(x)$ for all $x \in X$ iff the two functions $T \circ i : X \rightarrow W$ and $f : X \rightarrow W$ are equal.

1. *Diagram Completion Property.* For any F -vector space W and any function $f : X \rightarrow W$, there exists a unique F -linear map $T : V \rightarrow W$ such that the following diagram commutes.

$$\begin{array}{ccc} X & \xrightarrow{i} & V \\ & \searrow f & \downarrow T \\ & & W \end{array}$$

(Commutativity of the diagram means, by definition, that $f = T \circ i$.)

2. *Unique Factorization Property.* For any F -vector space W and any function $f : X \rightarrow W$, there exists a unique F -linear map $T : V \rightarrow W$ such that the factorization $f = T \circ i$ holds.
3. *Bijection between Collections of Functions.* For any F -vector space W , there is a bijection from the collection

$$A = \{F\text{-linear maps } T : V \rightarrow W\}$$

onto the collection

$$B = \{\text{arbitrary maps } f : X \rightarrow W\}.$$

The bijection sends $T \in A$ to $T \circ i \in B$. The inverse bijection sends $f \in B$ to the unique F -linear map $T : V \rightarrow W$ that extends f .

The first two restatements follow immediately from the original theorem and the fact that T extends f iff $f = T \circ i$ iff the given diagram commutes. To prove the third restatement, define a map ϕ with domain A by $\phi(T) = T \circ i$ for $T \in A$. For any such T , $T \circ i$ is a function from X to W , so that ϕ does map A into the codomain B . The unique factorization property amounts to the statement that $\phi : A \rightarrow B$ is a bijection. For, given any $f \in B$, that property says there exists a unique $T \in A$ with $f = T \circ i = \phi(T)$. Existence of such a T means that ϕ is onto; uniqueness of T means that ϕ is one-to-one. So, ϕ is a bijection. The relation $f = \phi(T)$ now implies that $\phi^{-1}(f) = T$, where T is the unique F -linear map extending f .

The theorem and its three restatements are all referred to as the *universal mapping property* for the basis X of the vector space V . More precisely, we should refer to this result as the UMP for the inclusion map $i : X \rightarrow V$. Note that composition with i establishes a canonical correspondence (bijection) between *F -linear* functions from V to another space W and *arbitrary* functions from X to W . The term “universal” indicates that the same map $i : X \rightarrow V$ works for all possible target vector spaces W .

The rest of this chapter discusses more examples of UMP’s that arise in linear and abstract algebra. Even more UMP’s are discussed in the following chapter on multilinear algebra. All of these universal mapping properties involve variations of the setup illustrated by the preceding example. In that example, we are interested in understanding special types of functions (namely, F -linear maps) from the fixed vector space V into arbitrary vector spaces W . The universal mapping property helps us understand these functions by setting up a bijection (for each fixed W) between these linear maps and another, simpler set of maps (namely, the set of all functions from the finite set X into W). The central idea of studying a linear map by computing its matrix relative to an ordered basis is really a manifestation of this universal mapping property (the columns of the matrix give the coordinates of $f(x_j)$ for each $x_j \in X$). Other UMP’s have a similar purpose: roughly speaking, “composition with a universal map” induces a bijection from collections of functions having one kind of structure to collections of functions having another kind of structure. The bijection is valuable since one of the two kinds of functions might be much easier to understand and analyze than the other. For instance, we will see UMP’s in the next chapter that convert “multilinear” maps (and their variations) into linear maps, thus permitting us to use linear algebra in the study of the more difficult subject of multilinear algebra.

19.1 Bases of Free R -Modules

The rest of this chapter assumes familiarity with the material on modules covered in Chapter 17. Our introductory example, involving bases of finite-dimensional vector spaces, readily generalizes to free modules over an arbitrary ring. Let R be a ring, and let M be a free left R -module with basis X . Recall this means that for each $v \in M$, there exists a unique expression $v = \sum_{x \in X} c_x x$ in which $c_x \in R$ and all but finitely many scalars c_x are zero. Let $i : X \rightarrow M$ be the inclusion map given by $i(x) = x$ for all $x \in X$. We have the following equivalent versions of the universal mapping property for $i : X \rightarrow M$.

1. *Diagram Completion Property.* For any left R -module N and any function $f : X \rightarrow N$, there exists a unique R -module homomorphism $T : M \rightarrow N$ such that the following diagram commutes.

$$\begin{array}{ccc} X & \xrightarrow{i} & M \\ & \searrow f & \downarrow T \\ & & N \end{array}$$

2. *Unique Factorization Property.* For any left R -module N and any function $f : X \rightarrow N$, there exists a unique R -module homomorphism $T : M \rightarrow N$ such that the factorization $f = T \circ i$ holds.
3. *Bijection between Collections of Functions.* For any left R -module N , there is a bijection from the collection

$$A = \{R\text{-module homomorphisms } T : M \rightarrow N\}$$

onto the collection

$$B = \{\text{arbitrary maps } f : X \rightarrow N\}.$$

The bijection sends $T \in A$ to $T \circ i \in B$. The inverse bijection sends $f \in B$ to the unique R -module homomorphism $T : M \rightarrow N$ that extends f .

The proof proceeds as before. Given N and $f : X \rightarrow N$, the only map $T : M \rightarrow N$ that could possibly be an R -module homomorphism extending f must be given by the formula

$$T \left(\sum_{x \in X} c_x x \right) = \sum_{x \in X} c_x f(x), \quad (c_x \in R),$$

which proves uniqueness of T . To prove existence of T , use the preceding formula as the definition of T . Then T is a well-defined function from M into N (the definition is unambiguous, since M is free with basis X), and one checks (as we did for vector spaces) that T is, indeed, an R -module homomorphism. This proves the first version of the UMP, and the other two versions follow as before.

19.2 Homomorphisms out of Quotient Modules

Let R be a ring, and let M be any left R -module with submodule N . Let $p : M \rightarrow M/N$ be the canonical R -homomorphism given by $p(x) = x + N$ for all $x \in M$. The following

universal mapping property of p helps explain the significance of the construction of the quotient module M/N . As before, we can state the UMP in several equivalent ways.

1. *UMP for Projection onto M/N (Diagram Completion Formulation):* For every left R -module Q and every R -linear map $f : M \rightarrow Q$ such that $f(x) = 0$ for all $x \in N$, there exists a unique R -module homomorphism $f' : M/N \rightarrow Q$ such that $f = f' \circ p$; i.e., such that $f'(x + N) = f(x)$ for all $x \in M$. (Furthermore, $\text{img}(f') = \text{img}(f)$ and $\ker(f') = \ker(f)/N$.)

$$\begin{array}{ccc} M & \xrightarrow{p} & M/N \\ & \searrow f & \downarrow f' \\ & & Q \end{array}$$

2. *UMP for Projection onto M/N (Bijective Formulation):* For any left R -module Q , there is a bijection from the collection

$$A = \{\text{all } R\text{-module homomorphisms } f' : M/N \rightarrow Q\}$$

onto the collection

$$B = \{R\text{-module homomorphisms } f : M \rightarrow Q \text{ such that } f(x) = 0 \text{ for all } x \in N\}.$$

The bijection sends $f' \in A$ to $f' \circ p \in B$. The inverse bijection sends $f \in B$ to the unique R -module homomorphism $f' : M/N \rightarrow Q$ such that $f' \circ p = f$.

To prove the diagram completion property, let $f : M \rightarrow Q$ satisfy $f(x) = 0$ for all $x \in N$. Uniqueness of $f' : M/N \rightarrow Q$ follows from the requirement that $f = f' \circ p$. For, this requirement means that $f'(x + N) = f'(p(x)) = (f' \circ p)(x) = f(x)$ for all $x \in M$, and every element of the set M/N has the form $x + N$ for some $x \in M$. Thus, if f' exists at all, it must be defined by the formula $f'(x + N) = f(x)$ for $x \in M$. Naturally, then, this is the formula we must use in the proof that f' does exist. To see that the formula defines a single-valued function, suppose $x, y \in M$ are such that $x + N = y + N$. Then $x - y \in N$, so that $f(x - y) = 0$ by assumption on f . Since f is a homomorphism, we have $f(x) - f(y) = f(x - y) = 0$ and hence $f(x) = f(y)$. Therefore, $f'(x + N) = f(x) = f(y) = f'(y + N)$, proving that $f' : M/N \rightarrow Q$ is a well-defined (single-valued) function. Knowing this, we can now check that f' is a R -homomorphism: for all $x, z \in M$ and $r \in R$, we have:

$$f'((x + N) + (z + N)) = f'((x + z) + N) = f(x + z) = f(x) + f(z) = f'(x + N) + f'(z + N);$$

$$f'(r(x + N)) = f'((rx) + N) = f(rx) = rf(x) = rf'(x + N).$$

As seen in the uniqueness proof, the very definition of f' guarantees that $f = f' \circ p$.

To verify the bijective version of the UMP, consider the function ϕ with domain A given by $\phi(f') = f' \circ p$ for $f' \in A$. Does ϕ map into the set B ? Yes, because: $f' \circ p$ is a R -homomorphism from M to Q , being a composition of R -homomorphisms; and this R -homomorphism does send every $x \in N$ to zero, since $(f' \circ p)(x) = f'(x + N) = f'(0 + N) = 0_Q$. The existence and uniqueness assertions proved above show that ϕ is a bijection.

Finally, we prove the parenthetical remark about the image and kernel of f' . We have

$$\text{img}(f') = \{f'(w) : w \in M/N\} = \{f'(x + N) : x \in M\} = \{f(x) : x \in M\} = \text{img}(f).$$

Next, note that $N \subseteq \ker(f)$ by our assumption on f , so the quotient module $\ker(f)/N$ makes sense. For $x \in M$, we have $x + N \in \ker(f')$ iff $f'(x + N) = 0$ iff $f(x) = 0$ iff $x \in \ker(f)$ iff $x + N \in \ker(f)/N$.

19.3 Direct Product of Two Modules

Let R be a ring, and let M and N be left R -modules. We can form the direct product $M \times N = M \oplus N = \{(m, n) : m \in M, n \in N\}$, which is a left R -module under componentwise operations (see §17.5). We define four canonical R -module homomorphisms associated with the module $M \times N$:

$$\begin{array}{lll} p : M \times N \rightarrow M & \text{given by} & p((m, n)) = m \quad (m \in M, n \in N); \\ q : M \times N \rightarrow N & \text{given by} & q((m, n)) = n \quad (m \in M, n \in N); \\ i : M \rightarrow M \oplus N & \text{given by} & i(m) = (m, 0) \quad (m \in M); \\ j : N \rightarrow M \oplus N & \text{given by} & j(n) = (0, n) \quad (n \in N). \end{array}$$

We call p and q the *natural projections*, while i and j are the *natural injections*. These maps satisfy the following identities:

$$p \circ i = \text{id}_M; \quad q \circ j = \text{id}_N; \quad q \circ i = 0; \quad p \circ j = 0; \quad i \circ p + j \circ q = \text{id}_{M \oplus N}.$$

In this section and the next, we discuss two different universal mapping properties for $M \times N = M \oplus N$, one involving the natural projections, and the other involving the natural injections.

1. *UMP for Projections (Diagram Completion Formulation)*: Suppose Q is any left R -module, and we are given two R -homomorphisms $f : Q \rightarrow M$ and $g : Q \rightarrow N$. There exists a unique R -homomorphism $h : Q \rightarrow M \times N$ such that $f = p \circ h$ and $g = q \circ h$, i.e., such that this diagram commutes:

$$\begin{array}{ccccc} & & M \times N & & \\ & \swarrow f & \downarrow h & \searrow g & \\ Q & & & & N \end{array}$$

2. *UMP for Projections (Bijective Formulation)*: For any left R -module Q , there is a bijection from the collection

$$A = \{\text{all } R\text{-module homomorphisms } h : Q \rightarrow M \times N\} = \text{Hom}_R(Q, M \times N)$$

onto the collection

$$\begin{aligned} B &= \{\text{pairs } (f, g) \text{ of } R\text{-homomorphisms } f : Q \rightarrow M, g : Q \rightarrow N\} \\ &= \text{Hom}_R(Q, M) \times \text{Hom}_R(Q, N). \end{aligned}$$

The bijection sends $h \in A$ to the pair $(p \circ h, q \circ h) \in B$. The inverse bijection sends $(f, g) \in B$ to the unique R -module homomorphism $h : Q \rightarrow M \times N$ such that $(f, g) = (p \circ h, q \circ h)$.

To verify the first version of the UMP, we start by proving the uniqueness of $h : Q \rightarrow M \times N$. If the function h exists at all, it must have the form $h(x) = (h_1(x), h_2(x))$ for all $x \in Q$, where $h_1 : Q \rightarrow M$ and $h_2 : Q \rightarrow N$ are certain functions. The requirement on h that $p \circ h = f$ forces $f(x) = p(h(x)) = p((h_1(x), h_2(x))) = h_1(x)$ for all $x \in Q$, so that $h_1 = f$. Similarly, the requirement $q \circ h = g$ forces $h_2 = g$. Therefore, h (if it exists at all) must be given by the formula $h(x) = (f(x), g(x))$ for $x \in Q$, so that uniqueness of h is proved.

To prove existence of h , define h (as we must) by setting $h(x) = (f(x), g(x))$ for $x \in Q$. By a computation similar to the one in the last paragraph, we have $p \circ h = f$ and $q \circ h = g$. We must still prove that h is an R -module homomorphism. For $x, y \in Q$ and $r \in R$, calculate:

$$\begin{aligned} h(x+y) &= (f(x+y), g(x+y)) = (f(x)+f(y), g(x)+g(y)) \\ &= (f(x), g(x)) + (f(y), g(y)) = h(x) + h(y); \\ h(rx) &= (f(rx), g(rx)) = (rf(x), rg(x)) = r(f(x), g(x)) = rh(x). \end{aligned}$$

To verify the bijective version of the UMP, consider the function ϕ with domain A given by $\phi(h) = (p \circ h, q \circ h)$ for $h \in A$. Does ϕ map into the set B ? Yes, because $p \circ h$ is an R -homomorphism from Q to M , and $q \circ h$ is an R -homomorphism from Q to N . The existence and uniqueness assertions proved above show that ϕ is a bijection.

19.4 Direct Sum of Two Modules

Keep the notation and assumptions from the previous section. We now consider the universal mapping property of the natural injections from M and N into $M \times N = M \oplus N$.

1. *UMP for Injections (Diagram Completion Formulation):* Suppose Q is any left R -module, and we have two R -homomorphisms $f : M \rightarrow Q$ and $g : N \rightarrow Q$. There exists a unique R -homomorphism $h : M \oplus N \rightarrow Q$ such that $f = h \circ i$ and $g = h \circ j$, i.e., such that this diagram commutes:

$$\begin{array}{ccccc} & & M & \xrightarrow{i} & M \oplus N & \xleftarrow{j} & N \\ & & \searrow f & & \downarrow h & & \swarrow g \\ & & & & Q & & \end{array}$$

2. *UMP for Injections (Bijective Formulation):* For any left R -module Q , there is a bijection from the collection

$$A = \{\text{all } R\text{-module homomorphisms } h : M \oplus N \rightarrow Q\} = \text{Hom}_R(M \oplus N, Q)$$

onto the collection

$$\begin{aligned} B &= \{\text{pairs } (f, g) \text{ of } R\text{-homomorphisms } f : M \rightarrow Q, g : N \rightarrow Q\} \\ &= \text{Hom}_R(M, Q) \times \text{Hom}_R(N, Q). \end{aligned}$$

The bijection sends $h \in A$ to the pair $(h \circ i, h \circ j) \in B$. The inverse bijection sends $(f, g) \in B$ to the unique R -module homomorphism $h : M \oplus N \rightarrow Q$ such that $(f, g) = (h \circ i, h \circ j)$.

To verify the first version of the UMP, we start by proving the uniqueness of $h : M \oplus N \rightarrow Q$. The requirement $h \circ i = f$ means that $h((m, 0)) = h(i(m)) = f(m)$ for all $m \in M$. The requirement $h \circ j = g$ means that $h((0, n)) = h(j(n)) = g(n)$ for all $n \in N$. Now, since h is also required to be an R -homomorphism, we must have (for all $m \in M$ and $n \in N$)

$$h((m, n)) = h((m, 0) + (0, n)) = h((m, 0)) + h((0, n)) = f(m) + g(n)$$

if h exists at all. This proves uniqueness of h .

To prove existence of $h : M \oplus N \rightarrow Q$, define h (as we must) by setting $h((m, n)) = f(m) + g(n)$ for all $m \in M$ and all $n \in N$. For each $m \in M$, we have $(h \circ i)(m) = h(i(m)) = h((m, 0)) = f(m) + g(0) = f(m) + 0 = f(m)$, so that $h \circ i = f$. Similarly, $f(0) = 0$ implies that $h \circ j = g$. Finally, is h an R -homomorphism? Let $m, m' \in M$, $n, n' \in N$, $r \in R$, and calculate:

$$\begin{aligned} h((m, n) + (m', n')) &= h((m + m', n + n')) = f(m + m') + g(n + n') \\ &= f(m) + f(m') + g(n) + g(n') = f(m) + g(n) + f(m') + g(n') \\ &= h((m, n)) + h((m', n')); \end{aligned} \quad (19.2)$$

$$h(r(m, n)) = h((rm, rn)) = f(rm) + g(rn) = rf(m) + rg(n) = r(f(m) + g(n)) = rh((m, n)).$$

Note that commutativity of addition in Q was needed to go from the first line to the second line in (19.2).

To verify the bijective version of the UMP, consider the function ϕ with domain A given by $\phi(h) = (h \circ i, h \circ j)$ for $h \in A$. Does ϕ map into the set B ? Yes, because $h \circ i$ is a R -homomorphism from M to Q , and $h \circ j$ is a R -homomorphism from N to Q . The existence and uniqueness assertions proved above show that ϕ is a bijection.

One may check (Exercises 9 and 10) that all the constructions and results in this section and the preceding one extend to products of the form $M_1 \times \cdots \times M_n$, where there are only finitely many factors. In the next two sections, we generalize the universal mapping properties even further, discussing arbitrary direct products and direct sums of R -modules. When there are infinitely many nonzero factors, the direct product $\prod_{i \in I} M_i$ is distinct from the direct sum $\bigoplus_{i \in I} M_i$, so that the two universal mapping properties (one for projections, one for injections) will involve different R -modules.

19.5 Direct Products of Arbitrary Families of R -Modules

Let R be a ring, let I be an index set, let M_i be a left R -module for each $i \in I$, and let $M = \prod_{i \in I} M_i$ be the direct product of the M_i 's. Recall from §17.5 that elements $x \in M$ are functions $x : I \rightarrow \bigcup_{i \in I} M_i$ such that $x(i) \in M_i$ for all $i \in I$. It is often convenient to think of these functions as “ I -tuples” $(x_i : i \in I)$, particularly when $I = \{1, 2, \dots, n\}$. Module operations in M are defined pointwise: for $x, y \in M$, we have $(x + y)(i) = x(i) + y(i) \in M_i$ and $(rx)(i) = r(x(i)) \in M_i$ for all $i \in I$.

For each $i \in I$, we have the *natural projection map* $p_i : M \rightarrow M_i$, which sends $x \in M$ to $x(i) \in M_i$. Each p_i is an R -homomorphism, since for $x, y \in M$ and $r \in R$,

$$\begin{aligned} p_i(x + y) &= (x + y)(i) = x(i) + y(i) = p_i(x) + p_i(y); \\ p_i(rx) &= (rx)(i) = r(x(i)) = rp_i(x). \end{aligned}$$

The family of projection maps $\{p_i : i \in I\}$ satisfies the following universal mapping property.

1. *UMP for Direct Products of Modules (Diagram Completion Formulation):* Suppose Q is any left R -module, and for each $i \in I$ we have an R -homomorphism $f_i : Q \rightarrow M_i$. There exists a unique R -homomorphism $f : Q \rightarrow \prod_{i \in I} M_i$ such

that $f_i = p_i \circ f$ for all $i \in I$, i.e., such that these diagrams commute for all $i \in I$:

$$\begin{array}{ccc} M_i & \xleftarrow{p_i} & \prod_{i \in I} M_i \\ & \swarrow f_i & \uparrow f \\ Q & & \end{array}$$

2. *UMP for Direct Products of Modules (Bijective Formulation):* For any left R -module Q , there is a bijection from the collection

$$A = \text{Hom}_R\left(Q, \prod_{i \in I} M_i\right) = \left\{ R\text{-module homomorphisms } f : Q \rightarrow \prod_{i \in I} M_i\right\}$$

onto the collection

$$B = \prod_{i \in I} \text{Hom}_R(Q, M_i) = \{I\text{-tuples } (f_i : i \in I) \text{ of } R\text{-homomorphisms } f_i : Q \rightarrow M_i\}.$$

The bijection sends $f \in A$ to $(p_i \circ f : i \in I) \in B$. The inverse bijection sends $(f_i : i \in I) \in B$ to the unique R -module homomorphism $f : Q \rightarrow M$ such that $f_i = p_i \circ f$ for all $i \in I$.

The proof of this UMP is similar to the one given earlier for $M \times N$. Consider the diagram completion property for fixed Q and fixed R -maps $f_i : Q \rightarrow M_i$ ($i \in I$). If $f : Q \rightarrow M$ exists at all, note that $f(x)$ is a function from I to $\bigcup_{i \in I} M_i$ for each $x \in Q$. Furthermore, the requirement that $f_i = p_i \circ f$ means that $f_i(x) = p_i(f(x))$ for all $x \in Q$. By definition of p_i , $p_i(f(x)) = f(x)(i)$ is the value of the function $f(x)$ at the point $i \in I$. Therefore, f is completely determined by the requirements $f_i = p_i \circ f$ for $i \in I$: we must have $f(x)(i) = f_i(x)$ for all $i \in I$ and all $x \in Q$; or, in I -tuple notation,

$$f(x) = (f_i(x) : i \in I) \quad \text{for all } x \in Q.$$

This proves uniqueness of f .

Proceeding to the existence proof, we must use the formula just written as the definition of $f : Q \rightarrow M$. Since $f(x)(i) = f_i(x) \in M_i$ for all $i \in I$, $f(x)$ is a well-defined element of $M = \prod_{i \in I} M_i$, and it is true that $f_i = p_i \circ f$ for all $i \in I$. We need only check that $f : Q \rightarrow M$ is an R -module homomorphism. Let $x, y \in Q$ and $r \in R$, and calculate:

$$\begin{aligned} f(x+y) &= (f_i(x+y) : i \in I) = (f_i(x) + f_i(y) : i \in I) \\ &= (f_i(x) : i \in I) + (f_i(y) : i \in I) = f(x) + f(y); \\ f(rx) &= (f_i(rx) : i \in I) = (rf_i(x) : i \in I) \\ &= r(f_i(x) : i \in I) = rf(x). \end{aligned}$$

For the second version of the UMP, define a function ϕ with domain A by setting $\phi(f) = (p_i \circ f : i \in I)$ for $f \in A$. For each $i \in I$, $p_i \circ f$ is, indeed, an R -homomorphism from Q to M_i . So ϕ does map into B , and the first version of the UMP shows that ϕ is a bijection.

19.6 Direct Sums of Arbitrary Families of R -Modules

Let R be a ring, let I be an index set, let M_i be a left R -module for each $i \in I$, and let $M = \bigoplus_{i \in I} M_i$ be the direct sum of the M_i 's. Recall that M is the submodule of the direct

product $\prod_{i \in I} M_i$ consisting of those functions $x : I \rightarrow \bigcup_{i \in I} M_i$ such that $x(i) \neq 0_{M_i}$ for only finitely many indices i . If the index set I is finite, then $\bigoplus_{i \in I} M_i = \prod_{i \in I} M_i$, but these modules are distinct when I is infinite and all M_i 's are nonzero. The module $\bigoplus_{i \in I} M_i$ is also referred to as the *coproduct* of the R -modules M_i .

For each $i \in I$, we have the *natural injection map* $j_i : M_i \rightarrow M$, which sends $m \in M_i$ to the function $x \in M$ such that $x(i) = m$ and $x(k) = 0_{M_k}$ for all $k \neq i$ in I . Informally, $j_i(m)$ is the I -tuple with m in position i and zeroes in all other positions. Using the Kronecker delta notation, defined by $\delta_{k,i} = 1_R$ if $k = i$ and $\delta_{k,i} = 0_R$ if $k \neq i$, we can write $j_i(m) = (\delta_{k,i}m : k \in I)$ for $m \in M_i$. Each j_i is an R -homomorphism, since for $m, n \in M_i$ and $r \in R$, we have:

$$j_i(m + n) = (\delta_{k,i}(m + n) : k \in I) = (\delta_{k,i}m : k \in I) + (\delta_{k,i}n : k \in I) = j_i(m) + j_i(n);$$

$$j_i(rm) = (\delta_{k,i}(rm) : k \in I) = r(\delta_{k,i}m : k \in I) = rj_i(m).$$

Next, we claim that any $x \in M = \bigoplus_{i \in I} M_i$ can be written as follows:

$$x = \sum_{i \in I} j_i(x(i)).$$

The sum on the right side is a sum of elements of M , in which all but finitely many summands $j_i(x(i))$ are zero, since all but finitely many of the elements $x(i)$ are zero. To verify the claim that the two functions x and $\sum_{i \in I} j_i(x(i))$ are equal, we evaluate each of them at an arbitrary $k \in I$:

$$\left[\sum_{i \in I} j_i(x(i)) \right] (k) = \sum_{i \in I} j_i(x(i))(k) = \sum_{i \in I} \delta_{k,i}x(i) = x(k).$$

Now we are ready to state the UMP satisfied by the family of injection maps $\{j_i : i \in I\}$.

1. *UMP for Direct Sums of Modules (Diagram Completion Formulation):* Suppose Q is any left R -module, and for each $i \in I$ we have an R -homomorphism $f_i : M_i \rightarrow Q$. There exists a unique R -homomorphism $f : \bigoplus_{i \in I} M_i \rightarrow Q$ such that $f_i = f \circ j_i$ for all i , i.e., such that these diagrams commute for all $i \in I$:

$$\begin{array}{ccc} M_i & \xrightarrow{j_i} & \bigoplus_{i \in I} M_i \\ & \searrow f_i & \downarrow f \\ & & Q \end{array}$$

2. *UMP for Direct Sums of Modules (Bijective Formulation):* For any left R -module Q , there is a bijection from the collection

$$A = \text{Hom}_R \left(\bigoplus_{i \in I} M_i, Q \right) = \left\{ \text{R-module homomorphisms } f : \bigoplus_{i \in I} M_i \rightarrow Q \right\}$$

onto the collection

$$B = \prod_{i \in I} \text{Hom}_R(M_i, Q) = \{I\text{-tuples } (f_i : i \in I) \text{ of } R\text{-homomorphisms } f_i : M_i \rightarrow Q\}.$$

The bijection sends $f \in A$ to $(f \circ j_i : i \in I) \in B$. The inverse bijection sends $(f_i : i \in I) \in B$ to the unique R -module homomorphism $f : \bigoplus_{i \in I} M_i \rightarrow Q$ such that $f_i = f \circ j_i$ for all $i \in I$.

We begin by proving uniqueness of f in the diagram completion version of the UMP. Fix Q and the R -maps $f_i : M_i \rightarrow Q$ for $i \in I$. Suppose $f : M = \bigoplus_{i \in I} M_i \rightarrow Q$ is any R -homomorphism such that $f \circ j_i = f_i$ for all i . Take any function $x \in M$, and write $x = \sum_{i \in I} j_i(x(i))$, as above. We must have

$$\begin{aligned} f(x) &= f\left(\sum_{i \in I} j_i(x(i))\right) = \sum_{i \in I} f(j_i(x(i))) \\ &= \sum_{i \in I} (f \circ j_i)(x(i)) = \sum_{i \in I} f_i(x(i)). \end{aligned}$$

This proves that f is uniquely determined on M by the f_i 's, if f exists at all. (Note that all sums written here and below are really finite sums, since we disregard all zero summands. None of the calculations in this proof make any sense for infinite direct products.)

To prove existence of f , define $f(x) = \sum_{i \in I} f_i(x(i))$ for all $x \in M$. This sum is a finite sum of elements of Q , since x is an element of the direct sum of the M_i 's, so that $f : M \rightarrow Q$ is a well-defined function. To confirm that $f \circ j_k = f_k$ for fixed $k \in I$, let us check that these functions agree at each $y \in M_k$. First, $(f \circ j_k)(y) = f(j_k(y)) = \sum_{i \in I} f_i(j_k(y)(i))$. Now, $j_k(y)(i) = 0$ if $i \neq k$, while $j_k(y)(k) = y$. Therefore, the sum has at most one nonzero summand, corresponding to $i = k$, and $(f \circ j_k)(y) = f_k(j_k(y)(k)) = f_k(y)$. It remains to check that $f : M \rightarrow Q$ is an R -module homomorphism. Let $x, y \in M$ and $r \in R$, and calculate:

$$\begin{aligned} f(x+y) &= \sum_{i \in I} f_i((x+y)(i)) = \sum_{i \in I} f_i(x(i)+y(i)) \\ &= \sum_{i \in I} [f_i(x(i)) + f_i(y(i))] = \sum_{i \in I} f_i(x(i)) + \sum_{i \in I} f_i(y(i)) = f(x) + f(y); \\ f(rx) &= \sum_{i \in I} f_i((rx)(i)) = \sum_{i \in I} f_i(r(x(i))) \\ &= \sum_{i \in I} r f_i(x(i)) = r \sum_{i \in I} f_i(x(i)) = rf(x). \end{aligned}$$

For the second version of the UMP, define a function ϕ with domain A by setting $\phi(f) = (f \circ j_i : i \in I)$ for $f \in A$. For each $i \in I$, $f \circ j_i$ is, indeed, an R -homomorphism from M_i to Q . So ϕ does map into B , and the first version of the UMP shows that ϕ is a bijection.

Now let I and K be index sets, and let M_i (for $i \in I$) and N_k (for $k \in K$) be left R -modules. By combining the bijections discussed in this section and the previous one, we obtain a bijection

$$\text{Hom}_R\left(\bigoplus_{i \in I} M_i, \prod_{k \in K} N_k\right) \rightarrow \prod_{i \in I} \prod_{k \in K} \text{Hom}_R(M_i, N_k)$$

that maps an R -module homomorphism $g : \bigoplus_{i \in I} M_i \rightarrow \prod_{k \in K} N_k$ to the tuple

$$(p_k \circ g \circ j_i : i \in I, k \in K),$$

where the p_k 's are the projections of $\prod_k N_k$ and the j_i 's are the injections of $\bigoplus_{i \in I} M_i$. Note that $p_k \circ g \circ j_i$ is an R -map from M_i to N_k . Given R -maps $g_{i,k} : M_i \rightarrow N_k$ for all $i \in I$ and $k \in K$, the inverse bijection maps the tuple $(g_{i,k} : i \in I, k \in K)$ to the R -homomorphism from $\bigoplus_{i \in I} M_i$ to $\prod_{k \in K} N_k$ such that

$$(x_i : i \in I) \mapsto \left(\sum_{i \in I} g_{i,k}(x_i) : k \in K \right).$$

This function can also be written

$$\left(x \mapsto \left(\sum_{i \in I} g_{i,k}(p'_i(x)) : k \in K \right) : x \in \bigoplus_{i \in I} M_i \right),$$

where p'_i denotes the projection of $\bigoplus_{j \in I} M_j$ onto M_i .

19.7 Solving Universal Mapping Problems

In this chapter, we have analyzed some basic algebraic constructions for modules and discovered the universal mapping properties (UMP's) of these constructions. The next chapter will adopt the opposite point of view: we start by specifying some *universal mapping problem* (also abbreviated UMP), and we then seek to construct a new object and map(s) that solve this problem. If we succeed in finding a solution to the UMP, it is also natural to inquire to what extent our solution is unique.

For example, our discussion of direct sums (coproducts) of R -modules suggests the analogous universal mapping problem for sets:

Problem (Coproducts for Sets). Given a family of sets $\{S_i : i \in I\}$, construct a set S and maps $j_i : S_i \rightarrow S$ satisfying the following UMP: for any set T and any collection of functions $g_i : S_i \rightarrow T$, there exists a unique function $g : S \rightarrow T$ with $g_i = g \circ j_i$ for all $i \in I$.

$$\begin{array}{ccc} S_i & \xrightarrow{j_i} & S \\ & \searrow g_i & \downarrow g \\ & & T \end{array}$$

Equivalently: construct a set S and maps $j_i : S_i \rightarrow S$ such that, for any set T , there is a bijection from the set

$$A = \{\text{all functions } g : S \rightarrow T\}$$

onto the set

$$B = \{\text{families of functions } (g_i : i \in I) \text{ where } g_i : S_i \rightarrow T\}$$

given by $g \mapsto (g \circ j_i : i \in I)$.

Note that we cannot form the “direct sum” of the S_i , since the S_i are sets, not R -modules. So a modification of our construction for R -modules is required. Here is one possible solution.

Construction of Solution to UMP. Let S be the disjoint union of the S_i ; formally, define

$$S = \{(i, x) : i \in I, x \in S_i\}.$$

For $i \in I$, define the function $j_i : S_i \rightarrow S$ by $j_i(x) = (i, x)$ for all $x \in S_i$. We must verify that S and the j_i 's do have the necessary universal mapping property. Assume T and $g_i : S_i \rightarrow T$ are given. To prove uniqueness, consider any function $g : S \rightarrow T$ satisfying $g_i = g \circ j_i$ for all $i \in I$. For $i \in I$ and $x \in S_i$, we then have $g((i, x)) = g(j_i(x)) = g_i(x)$. Thus, the value of g at every $(i, x) \in S$ is completely determined by the g_i 's. So g is unique if it exists at all.

To prove existence, we must define $g((i, x)) = g_i(x)$ for all $(i, x) \in S$. It is immediate that $g : S \rightarrow T$ is a well-defined function such that $g_i = g \circ j_i$ for all $i \in I$. The bijective version

of the UMP follows, as in earlier proofs, once we note that the function $g \mapsto (g \circ j_i : i \in I)$ does indeed map A into B .

Now we have “solved” the UMP posed above. But is our solution unique? Certainly not — we can always change notation to obtain superficially different solutions to the UMP. For instance, we could have defined $S = \{(x, i) : i \in I, x \in S_i\}$ and $j_i(x) = (x, i)$ for $i \in I$ and $x \in S_i$. On the other hand, we claim our solution is unique “up to a unique isomorphism compatible with the UMP.” In the case at hand, this means that for any other solution (S', j'_i) to the UMP, there exists a unique bijection $g : S \rightarrow S'$ such that $j'_i = g \circ j_i$ for all $i \in I$. Briefly, although ungrammatically, we say that the solution (S, j_i) to the UMP is *essentially unique*.

Proof of Essential Uniqueness of Solution. Suppose (S', j'_i) also solves the UMP. Applying the universal mapping property of (S, j_i) to the set $T = S'$ and the family of maps $g_i = j'_i$, we conclude at once that there exists a unique function $g : S \rightarrow S'$ with $j'_i = g \circ j_i$ for all $i \in I$.

$$\begin{array}{ccc} S_i & \xrightarrow{j_i} & S \\ & \searrow j'_i & \downarrow g \\ & & S' \end{array}$$

To complete the proof of essential uniqueness, we need only show that g is a bijection. For this, we can use the universal mapping property of (S', j'_i) to construct a candidate for the inverse of g . Specifically, let $T = S$ and $g_i = j_i$ in the UMP for (S', j'_i) . The UMP says that there exists a unique function $g' : S' \rightarrow S$ with $j_i = g' \circ j'_i$ for all $i \in I$.

$$\begin{array}{ccc} S_i & \xrightarrow{j'_i} & S' \\ & \searrow j_i & \downarrow g' \\ & & S \end{array}$$

Now, $\text{id}_S \circ j_i = j_i = g' \circ j'_i = (g' \circ g) \circ j_i$. Thus, $h = \text{id}_S$ and $h = g' \circ g$ are two functions from S to S with the property that $h \circ j_i = j_i$ for all $i \in I$. But according to the UMP for (S, j_i) (with $T = S$ and $g_i = j_i$), there is a *unique* map $h : S \rightarrow S$ with this property.

$$\begin{array}{ccc} S_i & \xrightarrow{j_i} & S \\ & \searrow j_i & \downarrow h \\ & & S \end{array}$$

Therefore, $g' \circ g = \text{id}_S$.

Similarly, $\text{id}_{S'} \circ j'_i = j'_i = g \circ j_i = (g \circ g') \circ j'_i$. So, $h = \text{id}_{S'}$ and $h = g \circ g'$ are two functions from S' to S' such that $h \circ j'_i = j'_i$ for all $i \in I$. But according to the UMP for (S', j'_i) (with $T = S'$ and $g_i = j'_i$), there is a *unique* map $h : S' \rightarrow S'$ with this property.

$$\begin{array}{ccc} S_i & \xrightarrow{j'_i} & S' \\ & \searrow j'_i & \downarrow h \\ & & S' \end{array}$$

Therefore, $g \circ g' = \text{id}_{S'}$. We now see that g' is the two-sided inverse of g , so both functions are bijections. This completes the proof.

The next chapter further develops the ideas presented here by posing and solving some universal mapping problems that appear at the foundations of multilinear algebra. For each UMP, we give an explicit construction showing that a solution to the UMP does exist. In each case, an argument completely analogous to the one just given proves that our solution is essentially unique, up to a unique isomorphism compatible with the universal maps. Once the universal mapping properties are available, we will use them to derive the basic facts about the algebraic structures occurring in multilinear algebra.

19.8 Summary

Here we summarize the universal mapping properties discussed in this chapter. We state each result as a diagram completion property and as a bijection between appropriate collections of functions.

1. *UMP for Basis of a Finite-Dimensional Vector Space.* Let $X = \{x_1, \dots, x_n\}$ be a basis of the vector space V over the field F , and let $i : X \rightarrow V$ be the inclusion map. For each F -vector space W , composition with i gives a bijection from the set $\text{Hom}_F(V, W)$ of all F -linear maps $T : V \rightarrow W$ onto the set of all functions $f : X \rightarrow W$. So, for each function $f : X \rightarrow W$ there exists a unique F -linear map $T : V \rightarrow W$ extending f ($f = T \circ i$):

$$\begin{array}{ccc} X & \xrightarrow{i} & V \\ & \searrow f & \downarrow T \\ & & W \end{array}$$

Explicitly, $T(\sum_{i=1}^n c_i x_i) = \sum_{i=1}^n c_i f(x_i)$ for $c_i \in F$.

2. *UMP for Basis of a Free Module.* Let R be a ring, let M be a free left R -module with basis X , and let $i : X \rightarrow M$ be the inclusion map. For each left R -module N , composition with i gives a bijection from the set $\text{Hom}_R(M, N)$ of all R -linear maps $T : M \rightarrow N$ onto the set of all functions $f : X \rightarrow N$. So, for each function $f : X \rightarrow N$ there exists a unique R -linear map $T : M \rightarrow N$ extending f ($f = T \circ i$):

$$\begin{array}{ccc} X & \xrightarrow{i} & M \\ & \searrow f & \downarrow T \\ & & N \end{array}$$

Explicitly, $T(\sum_{x \in X} c_x x) = \sum_{x \in X} c_x f(x)$ for $c_x \in R$ (where only finitely many c_x 's are nonzero).

3. *UMP for Quotient Modules.* Let R be a ring, let M be a left R -module, let N be a submodule of M , and let $p : M \rightarrow M/N$ be the canonical projection map. For each left R -module Q , composition with p gives a bijection from the set $\text{Hom}_R(M/N, Q)$ of all R -linear maps $f' : M/N \rightarrow Q$ onto the set of those R -linear maps $f : M \rightarrow Q$ satisfying $f(z) = 0$ for all $z \in N$. So, for each R -linear map f on M that sends all of N to zero, there exists a unique “lifting” of f to

an R -linear map f' on M/N ($f = f' \circ p$):

$$\begin{array}{ccc} M & \xrightarrow{p} & M/N \\ & \searrow f & \downarrow f' \\ & & Q \end{array}$$

Explicitly, $f'(x + N) = f(x)$ for all $x \in M$; moreover, $\text{img}(f') = \text{img}(f)$ and $\ker(f') = \ker(f)/N$.

4. *UMP for Direct Product of Two Modules.* Let R be a ring, let M and N be left R -modules, and let $p : M \times N \rightarrow M$ and $q : M \times N \rightarrow N$ be the natural projections. For each left R -module Q , the map $h \mapsto (p \circ h, q \circ h)$ is a bijection from $\text{Hom}_R(Q, M \times N)$ onto $\text{Hom}_R(Q, M) \times \text{Hom}_R(Q, N)$. So, for each pair of R -linear maps (f, g) with $f : Q \rightarrow M$ and $g : Q \rightarrow N$, there exists a unique R -linear $h : Q \rightarrow M \times N$ with $f = p \circ h$ and $g = q \circ h$:

$$\begin{array}{ccccc} M & \xleftarrow{p} & M \times N & \xrightarrow{q} & N \\ & \swarrow f & \uparrow h & \nearrow g & \\ & & Q & & \end{array}$$

Explicitly, $h(x) = (f(x), g(x))$ for $x \in Q$.

5. *UMP for Direct Sum of Two Modules.* Let R be a ring, let M and N be left R -modules, and let $i : M \rightarrow M \oplus N$ and $j : N \rightarrow M \oplus N$ be the natural injections. For each left R -module Q , the map $h \mapsto (h \circ i, h \circ j)$ is a bijection from $\text{Hom}_R(M \oplus N, Q)$ onto $\text{Hom}_R(M, Q) \times \text{Hom}_R(N, Q)$. So, for each pair of R -linear maps (f, g) with $f : M \rightarrow Q$ and $g : N \rightarrow Q$, there exists a unique R -linear $h : M \oplus N \rightarrow Q$ with $f = h \circ i$ and $g = h \circ j$:

$$\begin{array}{ccc} M & \xrightarrow{i} & M \oplus N & \xleftarrow{j} & N \\ & \searrow f & \downarrow h & \nearrow g & \\ & & Q & & \end{array}$$

Explicitly, $h((m, n)) = f(m) + g(n)$ for $m \in M$ and $n \in N$.

6. *UMP for Direct Product of a Family of Modules.* Let R be a ring, let I be an index set, let M_i be a left R -module for each $i \in I$, and let $p_i : \prod_{j \in I} M_j \rightarrow M_i$ be the natural projections. For each left R -module Q , the map $f \mapsto (p_i \circ f : i \in I)$ is a bijection from $\text{Hom}_R(Q, \prod_{j \in I} M_j)$ onto $\prod_{j \in I} \text{Hom}_R(Q, M_j)$. So, for each family of R -linear maps $(f_i : i \in I)$ with $f_i : Q \rightarrow M_i$ for all $i \in I$, there exists a unique R -linear $f : Q \rightarrow \prod_{j \in I} M_j$ with $f_i = p_i \circ f$ for all $i \in I$:

$$\begin{array}{ccc} M_i & \xleftarrow{p_i} & \prod_{j \in I} M_j \\ & \swarrow f_i & \uparrow f \\ & & Q \end{array}$$

Explicitly, $f(x) = (f_i(x) : i \in I)$ for $x \in Q$.

7. *UMP for Direct Sum of a Family of Modules.* Let R be a ring, let I be an index set, let M_i be a left R -module for each $i \in I$, and let $j_i : M_i \rightarrow \bigoplus_{k \in I} M_k$ be the natural injections. For each left R -module Q , the map $f \mapsto (f \circ j_i : i \in I)$ is a bijection from $\text{Hom}_R(\bigoplus_{k \in I} M_k, Q)$ onto $\prod_{k \in I} \text{Hom}_R(M_k, Q)$. So, for each family of R -linear maps $(f_i : i \in I)$ with $f_i : M_i \rightarrow Q$ for all $i \in I$, there exists a unique R -linear $f : \bigoplus_{k \in I} M_k \rightarrow Q$ with $f_i = f \circ j_i$ for all $i \in I$:

$$\begin{array}{ccc} M_i & \xrightarrow{j_i} & \bigoplus_{k \in I} M_k \\ & \searrow f_i & \downarrow f \\ & & Q \end{array}$$

Explicitly, $f(x) = \sum_{i \in I} f_i(x_i)$ for $x = (x_i : i \in I) \in \bigoplus_{k \in I} M_k$.

8. Combining the last two items, we have a bijection

$$\text{Hom}_R\left(\bigoplus_{i \in I} M_i, \prod_{k \in K} N_k\right) \longrightarrow \prod_{i \in I} \prod_{k \in K} \text{Hom}_R(M_i, N_k)$$

that sends an R -map $g : \bigoplus_{i \in I} M_i \rightarrow \prod_{k \in K} N_k$ to $(p_k \circ g \circ j_i : i \in I, k \in K)$.

9. *UMP for Coproduct of Sets.* Let I be an index set, and let S_i be a set for each $i \in I$. Define $S = \{(i, x) : i \in I, x \in S_i\}$ and define $j_i : S_i \rightarrow S$ by $j_i(x) = (i, x)$ for $i \in I$ and $x \in S_i$. For each set T , the map $f \mapsto (f \circ j_i : i \in I)$ is a bijection from the set of functions from S to T to the set of families $(f_i : i \in I)$ where $f_i : S_i \rightarrow T$ for $i \in I$. So, for each family $(f_i : i \in I)$ of functions from S_i to T , there exists a unique $f : S \rightarrow T$ with $f_i = f \circ j_i$ for all $i \in I$:

$$\begin{array}{ccc} S_i & \xrightarrow{j_i} & S \\ & \searrow f_i & \downarrow f \\ & & T \end{array}$$

Explicitly, $f((i, x)) = f_i(x)$ for $i \in I$ and $x \in S_i$.

The solution of a universal mapping problem is essentially unique, meaning that for any two solutions to a UMP, there is a unique bijection between them that respects the associated universal maps.

19.9 Exercises

In these exercises, assume R is a ring unless otherwise stated.

- Let V and W be nonzero finite-dimensional vector spaces over a field F . Let $X = \{x_1, \dots, x_n\}$ be a subset of V . Assume X spans V but is linearly dependent over F . (a) Prove or disprove: for every function $f : X \rightarrow W$, there exists an F -linear map $T : V \rightarrow W$ extending f . (b) Prove or disprove: for all $f : X \rightarrow W$, there is at most one F -linear map $T : V \rightarrow W$ extending f .
- Repeat Exercise 1, but now assume X is a linearly independent subset of V that does not span V .

3. Give complete details of the proof of the UMP for bases of a free R -module stated in §19.1. (Do not assume the basis X is finite.)
4. *UMP for Quotient Groups.* Let (G, \star) be a group with normal subgroup N . Let $p : G \rightarrow G/N$ be the projection $p(x) = x \star N$ for $x \in G$. (a) Prove: for every group L and every group homomorphism $f : G \rightarrow L$ such that $f(x) = e_L$ for all $x \in N$, there exists a unique group homomorphism $f' : G/N \rightarrow L$ such that $f = f' \circ p$. What are $\text{img}(f')$ and $\ker(f')$? (b) Restate (a) in terms of a bijection between two collections of functions.
5. *UMP for Quotient Rings.* Let $(R, +, \cdot)$ be a ring with ideal I . Let $p : R \rightarrow R/I$ be the projection $p(x) = x + I$ for $x \in R$. Formulate and prove a universal mapping property characterizing the quotient ring R/I and the map $p : R \rightarrow R/I$.
6. Let M and N be left R -modules. Check carefully that the natural projections and natural injections for $M \times N = M \oplus N$ are R -linear and satisfy the following identities:

$$p \circ i = \text{id}_M; \quad q \circ j = \text{id}_N; \quad q \circ i = 0; \quad p \circ j = 0; \quad i \circ p + j \circ q = \text{id}_{M \oplus N}.$$

7. Assume R is a commutative ring. (a) Prove that the bijection

$$\phi : \text{Hom}_R(Q, M \times N) \rightarrow \text{Hom}_R(Q, M) \times \text{Hom}_R(Q, N)$$

constructed in §19.3 is an R -linear map. (b) Prove that $\text{Hom}_R(Q, \prod_{i \in I} M_i) \cong \prod_{i \in I} \text{Hom}_R(Q, M_i)$ as R -modules, via the bijection in §19.5.

8. Assume R is a commutative ring. (a) Prove that the bijection

$$\phi : \text{Hom}_R(M \times N, Q) \rightarrow \text{Hom}_R(M, Q) \times \text{Hom}_R(N, Q)$$

constructed in §19.4 is an R -linear map. (b) Prove that $\text{Hom}_R(\bigoplus_{i \in I} M_i, Q) \cong \prod_{i \in I} \text{Hom}_R(M_i, Q)$ as R -modules, via the bijection in §19.6.

9. For fixed $k \in \mathbb{N}^+$, let M_1, M_2, \dots, M_k be left R -modules. Prove that for all left R -modules Q , there is a bijection from $\text{Hom}_R(Q, M_1 \times M_2 \times \dots \times M_k)$ to $\text{Hom}_R(Q, M_1) \times \text{Hom}_R(Q, M_2) \times \dots \times \text{Hom}_R(Q, M_k)$: (a) by using induction on k and the bijections in §19.3; (b) by imitating the construction in §19.3.
10. For fixed $k \in \mathbb{N}^+$, let M_1, M_2, \dots, M_k be left R -modules. Prove that for all left R -modules Q , there is a bijection from $\text{Hom}_R(M_1 \oplus M_2 \oplus \dots \oplus M_k, Q)$ to $\text{Hom}_R(M_1, Q) \times \text{Hom}_R(M_2, Q) \times \dots \times \text{Hom}_R(M_k, Q)$: (a) by using induction on k and the bijections in §19.4; (b) by imitating the construction in §19.4.
11. Assume R is a commutative ring. With the setup in §19.1, show that the bijection $T \mapsto T \circ i$ is, in fact, an R -module isomorphism between the R -module $\text{Hom}_R(M, N)$ and the product R -module N^X .
12. With the setup in §19.5, let $f : Q \rightarrow \prod_{i \in I} M_i$ be the R -map corresponding to a given family of R -maps $f_i : Q \rightarrow M_i$. Prove $\ker(f) = \bigcap_{i \in I} \ker(f_i)$.
13. With the setup in §19.6, let $f : \bigoplus_{i \in I} M_i \rightarrow Q$ be the R -map corresponding to a given family of R -maps $f_i : M_i \rightarrow Q$. Prove $\text{img}(f) = \sum_{i \in I} \text{img}(f_i)$.
14. *UMP for Products of Sets.* Given an index set I and a set S_i for each $i \in I$, formulate and prove a universal mapping property satisfied by the Cartesian product set $S = \prod_{i \in I} S_i$ and the natural projection functions $p_i : S \rightarrow S_i$ given by $p_i((x_k : k \in I)) = x_i$ for each $i \in I$.

15. *UMP for Products of Groups and Rings.* (a) Given a family of groups $\{G_i : i \in I\}$ (not necessarily commutative), construct a group G and group homomorphisms $p_i : G \rightarrow G_i$ (for all $i \in I$) such that, for any group K , there is a bijection from the set A of all group homomorphisms $f : K \rightarrow G$ onto the set B of all families of group homomorphisms $(f_i : i \in I)$ with $f_i : K \rightarrow G_i$ for all $i \in I$, given by $f \mapsto (p_i \circ f : i \in I)$ for $f \in A$. (b) Does the construction in (a) still work if we replace groups and group homomorphisms by rings and ring homomorphisms throughout?
16. In §19.4, we showed that for any two left R -modules M and N , the direct sum $M \oplus N$ and the natural injections $i : M \rightarrow M \oplus N$ and $j : N \rightarrow M \oplus N$ satisfy the UMP for the coproduct of two modules. (a) Does this construction of the coproduct still work if we assume M and N are groups (possibly non-commutative), and demand that all maps be group homomorphisms? (b) Does this construction of the coproduct still work if we assume M and N are rings and demand that all maps be ring homomorphisms? What if we only allow commutative rings?
17. Let M be a left R -module with submodule N , and suppose there is a left R -module Z and a map $q : M \rightarrow Z$ such that $q[N] = \{0_Z\}$, and Z and q satisfy the UMP for quotient modules from §19.2. Prove that there exists a unique R -module isomorphism $g : M/N \rightarrow Z$ with $q = g \circ p$, where $p : M \rightarrow M/N$ is the natural projection. (Imitate the essential uniqueness proof in §19.7.)
18. Give a specific example using \mathbb{Z} -modules to show that the result of Exercise 17 might fail without the hypothesis that $q[N] = \{0_Z\}$.
19. Let $\{M_i : i \in I\}$ be an indexed family of left R -modules. (a) Carefully state what it means to say that the direct product $\prod_{i \in I} M_i$ and the associated projection maps are *essentially unique*, and then prove it. (b) Carefully state what it means to say that the direct sum $\bigoplus_{i \in I} M_i$ and the associated injection maps are *essentially unique*, and then prove it.
20. Given left R -modules M and N , prove that $M \times N \cong N \times M$ using the fact that both modules solve the same universal mapping problem. Then find a formula for the unique isomorphism compatible with the natural injections.
21. We are given an index set I , left R -modules M_i and N_i for each $i \in I$, and an R -map $f_i : M_i \rightarrow N_i$ for each $i \in I$. Let $p_i : \prod_{k \in I} M_k \rightarrow M_i$ and $q_i : \prod_{k \in I} N_k \rightarrow N_i$ be the natural projection maps, for each $i \in I$. (a) Use the UMP for direct products to show there exists a unique R -map $F : \prod_{k \in I} M_k \rightarrow \prod_{k \in I} N_k$ such that $q_i \circ F = f_i \circ p_i$ for all $i \in I$. (b) Find an explicit formula for F , and describe $\ker(F)$ and $\text{img}(F)$. (c) Write $F = F(f_i : i \in I) = F(f_i)$ to indicate the dependence of F on the given maps f_i . Suppose P_i is a left R -module and $g_i : N_i \rightarrow P_i$ is an R -map, for all $i \in I$. Let $r_i : \prod_{k \in I} P_k \rightarrow P_i$ be the natural projection maps. Prove that $F(g_i \circ f_i) = F(g_i) \circ F(f_i)$ in two ways: using the explicit formula in (b), and using the uniqueness of F proved in (a).
22. With the setup in Exercise 21, state and prove results analogous to (a), (b), and (c) of that exercise for a map $G : \bigoplus_{k \in I} M_k \rightarrow \bigoplus_{k \in I} N_k$ induced by the f_i 's and compatible with the natural injections.
23. Let M , N , and P be free left R -modules with bases X , Y , and Z and inclusion maps $i : X \rightarrow M$, $j : Y \rightarrow N$, and $k : Z \rightarrow P$. (a) Show that for each function $f : X \rightarrow Y$, there exists a unique R -map $F(f) : M \rightarrow N$ with $F(f) \circ i = j \circ f$. Do not find a formula for $F(f)$. (b) Use (a) (and similar results for functions

- from Y to Z , etc.) to show that for all functions $f : X \rightarrow Y$ and $g : Y \rightarrow Z$, $F(g \circ f) = F(g) \circ F(f)$. (c) Use (a) to show that $F(\text{id}_X) = \text{id}_M$. [In the language of category theory, F is a *functor* from the category of sets and functions to the category of left R -modules and R -maps.]
24. *UMP for Quotient Sets.* Let X be a set, let \sim be an equivalence relation on X , and let X/\sim be the set of all equivalence classes of \sim . Call a function f with domain X *compatible with \sim* iff for all $x, y \in X$, $x \sim y$ implies $f(x) = f(y)$. Let $p : X \rightarrow X/\sim$ be the map that sends each $x \in X$ to its equivalence class $[x]$ relative to \sim , given by $[x] = \{z \in X : x \sim z\}$. Note that p is surjective and compatible with \sim . (a) Prove that for every set Z , the map $h \mapsto h \circ p$ defines a bijection from the set A of all functions $h : X/\sim \rightarrow Z$ to the set B of all functions $f : X \rightarrow Z$ that are compatible with \sim . (b) Prove that X/\sim (and the map p , compatible with \sim) is the essentially unique solution to the UMP in (a). (c) Explain how quotient modules are a special case of the construction in this problem (cf. §19.2).
25. *UMP for Quotient Topologies.* A *topological space* is a set X and a family of subsets of X , called *open sets*, such that: \emptyset and X are open; the union of any collection of open sets is open; and the intersection of finitely many open sets is open. A function $f : X \rightarrow Y$ between two such spaces is *continuous* iff for all open subsets V of Y , $f^{-1}[V]$ is an open subset of X . Let X be a topological space, and let \sim be an equivalence relation on X . Define X/\sim and $p : X \rightarrow X/\sim$ as in Exercise 24. Define $V \subseteq X/\sim$ to be an open set iff $p^{-1}[V] \subseteq X$ is open in the given topological space X . (a) Show that this definition makes X/\sim into a topological space and p into a continuous surjective map. (b) Prove that for every topological space Z , the map $h \mapsto h \circ p$ defines a bijection from the set of all *continuous* functions $h : X/\sim \rightarrow Z$ to the set of all *continuous* functions $f : X \rightarrow Z$ that are compatible with \sim . (Explain why it suffices, using Exercise 24(a), to show that h is continuous if and only if $h \circ p$ is continuous.)
26. *Products and Coproducts of Topological Spaces.* Let $\{X_i : i \in I\}$ be a family of topological spaces. (a) Construct a “product” of the spaces X_i satisfying a universal mapping property analogous to the product of modules (where all functions under consideration are required to be continuous). (b) Construct a “coproduct” of the spaces X_i satisfying a universal mapping property analogous to direct sums of modules and coproducts of sets (again, all functions considered must be continuous).
27. *UMP for Polynomial Rings.* (a) Let R be a commutative ring, and let $i : R \rightarrow R[x]$ be given by $i(r) = (r, 0, 0, \dots)$ for $r \in R$ (so i maps r to the constant polynomial with constant term r). Prove: for each commutative ring S , the map $H \mapsto (H \circ i, H(x))$ defines a bijection from the set of ring homomorphisms $H : R[x] \rightarrow S$ to the set of pairs (h, c) , where $h : R \rightarrow S$ is a ring homomorphism and $c \in S$ (cf. §3.5). (b) For $m \in \mathbb{N}^+$, generalize (a) to obtain a bijective formulation of the UMP for the polynomial ring $R[x_1, \dots, x_m]$ (cf. §3.20).
28. *Localization of a Commutative Ring.* For any commutative ring T , we say $t \in T$ is a *unit of T* iff there exists $u \in T$ with $tu = 1_T = ut$; let T^* be the set of units of T . The goal of this exercise is to solve the following universal mapping problem. Let R be a commutative ring, and let S be a subset of R that contains 1_R and is closed under the multiplication of R . Construct a commutative ring L and a ring homomorphism $i : R \rightarrow L$ such that $i(s) \in L^*$ for all $s \in S$; and for all commutative rings T , the map $g \mapsto g \circ i$ defines a bijection from the set of all ring

homomorphisms $g : L \rightarrow T$ onto the set of all ring homomorphisms $f : R \rightarrow T$ such that $f(s) \in T^*$ for all $s \in S$.

$$\begin{array}{ccc} R & \xrightarrow{i} & L \\ & \searrow f & \downarrow g \\ & T & \end{array} \quad \begin{array}{l} (i[S] \subseteq L^*) \\ (f[S] \subseteq T^*) \end{array}$$

Intuitively, we are trying to build a ring L in which all elements of S become invertible, and this ring should be “as close to R as possible.” The idea will be to use “fractions” r/s , with $r \in R$ and $s \in S$, as the elements of L . (a) Let $X = R \times S = \{(r, s) : r \in R, s \in S\}$. Define a binary relation \sim on X by setting $(r, s) \sim (r', s')$ iff there exists $t \in S$ with $t(rs' - sr') = 0$. Check that \sim is an equivalence relation on X . (b) Let L be the set of equivalence classes of \sim on X , and write r/s for the equivalence class of $(r, s) \in X$. For r/s and u/v in L , define $r/s + u/v = (rv + su)/(sv)$ and $(r/s) \cdot (u/v) = (ru)/(sv)$. Check that these two operations are well-defined. (c) Verify that $(L, +, \cdot)$ is a commutative ring. (d) Define $i : R \rightarrow L$ by setting $i(r) = r/1_R$ for all $r \in R$. Check that i is a ring homomorphism, and $i(s) \in L^*$ for all $s \in S$. (e) Verify that L and i solve the UMP described above. (f) Explain why L and i are essentially unique. (g) Explain why the construction of \mathbb{Q} from \mathbb{Z} is a special case of the construction in this exercise.

29. Define an *inverse chain* to be a collection of left R -modules $(M_n : n \in \mathbb{N})$ and R -maps $f_n : M_n \rightarrow M_{n-1}$ for $n \in \mathbb{N}^+$:

$$M_0 \xleftarrow{f_1} M_1 \xleftarrow{f_2} M_2 \xleftarrow{f_3} \cdots \xleftarrow{f_{n-1}} M_{n-1} \xleftarrow{f_n} M_n \xleftarrow{f_{n+1}} \cdots .$$

Solve the following universal mapping problem: construct a left R -module L and R -maps $p_n : L \rightarrow M_n$ for $n \in \mathbb{N}$ with $f_n \circ p_n = p_{n-1}$ for all $n \in \mathbb{N}^+$, such that for any left R -module P and R -maps $g_n : P \rightarrow M_n$ ($n \in \mathbb{N}$) with $f_n \circ g_n = g_{n-1}$ for all $n \in \mathbb{N}^+$, there exists a unique R -map $h : P \rightarrow L$ with $p_n \circ h = g_n$ for all $n \in \mathbb{N}$. L is called the *inverse limit* of the inverse chain. [Hint: Define L to be a certain submodule of $\prod_{n \in \mathbb{N}} M_n$.]

30. Define a *direct chain* to be a collection of left R -modules $(M_n : n \in \mathbb{N})$ and R -maps $f_n : M_n \rightarrow M_{n+1}$ for $n \in \mathbb{N}$:

$$M_0 \xrightarrow{f_0} M_1 \xrightarrow{f_1} M_2 \xrightarrow{f_2} \cdots \xrightarrow{f_{n-1}} M_n \xrightarrow{f_n} M_{n+1} \xrightarrow{f_{n+1}} \cdots .$$

Solve the following universal mapping problem: construct a left R -module D and R -maps $j_n : M_n \rightarrow D$ for $n \in \mathbb{N}$ with $j_{n+1} \circ f_n = j_n$ for all $n \in \mathbb{N}$, such that for any left R -module P and R -maps $g_n : M_n \rightarrow P$ ($n \in \mathbb{N}$) with $g_{n+1} \circ f_n = g_n$ for all $n \in \mathbb{N}$, there exists a unique R -map $h : D \rightarrow P$ with $h \circ j_n = g_n$ for all $n \in \mathbb{N}$. D is called the *direct limit* of the direct chain. [Hint: Define D to be a certain quotient module of $\bigoplus_{n \in \mathbb{N}} M_n$. Your proof will be cleaner if you use UMP’s for direct sums and quotient modules properly.]

31. *Free Monoid Generated by a Set.* A *monoid* is a pair (M, \star) satisfying the first three group axioms in Table 1.1 (closure, associativity, and identity). Given monoids M and N , a *monoid homomorphism* is a map $f : M \rightarrow N$ such that $f(xy) = f(x)f(y)$ for all $x, y \in M$, and $f(1_M) = 1_N$. Given any set X , our goal is to solve the following universal mapping problem (cf. §19.1): construct a monoid M and a function $i : X \rightarrow M$ such that for any monoid N and any function

$f : X \rightarrow N$, there exists a unique monoid homomorphism $T : M \rightarrow N$ with $f = T \circ i$. (a) Let M be the set of all finite sequences $w = w_1 w_2 \cdots w_k$ (called *words*) where $k \in \mathbb{N}$ and each $w_i \in X$. Note that the empty sequence, denoted ϵ , is in M . Given $w = w_1 w_2 \cdots w_k$ and $y = y_1 y_2 \cdots y_m$ in M , define $w \star y$ to be the concatenation $w_1 w_2 \cdots w_k y_1 y_2 \cdots y_m$. Show that (M, \star) is a monoid with identity ϵ . (b) Define $i : X \rightarrow M$ by letting $i(x)$ be the word (sequence) of length 1 with sole entry x , for $x \in X$. Prove that M and i solve the universal mapping problem posed above.

32. *Coproduct of Two Groups.* Given two groups M and N , we wish to construct a new group $M * N$ and group homomorphisms $i : M \rightarrow M * N$ and $j : N \rightarrow M * N$ satisfying an analogue of the universal mapping problem in §19.4 (with modules and R -maps replaced by groups and group homomorphisms). Let $M' = M \sim \{e_M\}$ be the set of non-identity elements in M ; similarly, let $N' = N \sim \{e_N\}$. If needed, change notation so that M' and N' are disjoint sets. Let the set $M * N$ consist of all “words” $w_1 w_2 \cdots w_k$ with $k \in \mathbb{N}$, $w_i \in M' \cup N'$ for all i , and (for $i < k$) $w_i \in M'$ iff $w_{i+1} \in N'$ (so letters alternate between M' and N'). Given $w = w_1 w_2 \cdots w_k$ and $y = y_1 y_2 \cdots y_m$ in $M * N$, define $w \star y$ by the following recursive rules. If $w = \epsilon$ (the empty word), then $w \star y = y$ (similarly if $y = \epsilon$). If $w_k \in M'$ and $y_1 \in N'$, or if $w_k \in N'$ and $y_1 \in M'$, then $w \star y$ is the concatenation $w_1 w_2 \cdots w_k y_1 y_2 \cdots y_m$. If $w_k, y_1 \in M'$ and $w_k y_1 = z \neq e_M$, then $w \star y = w_1 \cdots w_{k-1} z y_2 \cdots y_m$. If $w_k, y_1 \in M'$ and $w_k y_1 = e_M$, we recursively define $w \star y = (w_1 \cdots w_{k-1}) \star (y_2 \cdots y_m)$. If $w_k, y_1 \in N'$ and $w_k y_1 = z \neq e_N$, then $w \star y = w_1 \cdots w_{k-1} z y_2 \cdots y_m$. If $w_k, y_1 \in N'$ and $w_k y_1 = e_N$, we recursively define $w \star y = (w_1 \cdots w_{k-1}) \star (y_2 \cdots y_m)$. (a) Prove $(M * N, \star)$ is a group. [The verification of associativity is rather tricky.] (b) Define injective group homomorphisms $i : M \rightarrow M * N$ and $j : N \rightarrow M * N$, and prove that $M * N$ with these maps solves the UMP posed above.
33. Generalize the construction in the previous exercise by defining a coproduct of a family of groups $\{G_i : i \in I\}$ satisfying a UMP analogous to the one in §19.6.
34. *Free Group Generated by a Set X .* Given any set X , construct a group (G, \star) and a function $i : X \rightarrow G$ such that for any group K and any function $f : X \rightarrow K$, there exists a unique group homomorphism $T : G \rightarrow K$ with $f = T \circ i$. (Ideas from the previous three exercises can help here. Let G consist of certain words in the “alphabet” $X \cup X'$, where X' is a set with $|X'| = |X|$ and $X \cap X' = \emptyset$; the elements of X' represent formal inverses of elements of X .)

Universal Mapping Problems in Multilinear Algebra

This chapter gives an introduction to the subject of *multilinear algebra*, which studies functions of several variables that are linear in each variable. After defining these “multilinear” maps, as well as alternating maps and symmetric maps, we pose and solve universal mapping problems that convert these less familiar maps to linear maps. The modules that arise in these constructions are called tensor products, exterior powers, and symmetric powers.

After building the tensor product, we use its universal mapping property to prove some isomorphisms and other general facts about tensor products. We show that linear maps between modules induce associated linear maps on tensor products, exterior powers, and symmetric powers. In the case of free modules, we find bases for these new modules in terms of bases for the original modules. This leads to a discussion of tensor products of matrices and the relation between determinants and exterior powers. The chapter ends with the construction of the tensor algebra of a module. Throughout the whole development, we stress the use of universal mapping properties as a means of organizing, motivating, and proving the basic results of multilinear algebra.

To read this chapter, one should be familiar with facts about modules and universal mapping properties covered in Chapters 17 and 19, as well as properties of permutations from Chapter 2.

20.1 Multilinear Maps

Throughout this chapter, R will always denote a commutative ring. Given R -modules M_1, \dots, M_n , consider the product R -module $M = M_1 \times \dots \times M_n$. We want to distinguish two special types of maps from M to another R -module N . On one hand, recall that $f : M \rightarrow N$ is an *R -module homomorphism* or an *R -linear map* iff $f(m + m') = f(m) + f(m')$ and $f(rm) = rf(m)$ for all $m, m' \in M$ and all $r \in R$. Writing this out in terms of components, this means that for all $m_k, m'_k \in M_k$ and all $r \in R$,

$$f(m_1 + m'_1, m_2 + m'_2, \dots, m_n + m'_n) = f(m_1, m_2, \dots, m_n) + f(m'_1, m'_2, \dots, m'_n),$$

$$f(rm_1, rm_2, \dots, rm_n) = rf(m_1, m_2, \dots, m_n).$$

On the other hand, we now define $f : M_1 \times \dots \times M_n \rightarrow N$ to be *R -multilinear* iff f is R -linear in each of its n arguments separately. More precisely, for each i between 1 and n and each fixed choice of $m_1, \dots, m_{i-1}, m_{i+1}, \dots, m_n$, we require that

$$f(m_1, \dots, m_{i-1}, m_i + m'_i, m_{i+1}, \dots, m_n) = f(m_1, \dots, m_i, \dots, m_n) + f(m_1, \dots, m'_i, \dots, m_n), \quad (20.1)$$

$$f(m_1, \dots, m_{i-1}, rm_i, m_{i+1}, \dots, m_n) = rf(m_1, \dots, m_i, \dots, m_n) \quad (20.2)$$

for all $m_i, m'_i \in M_i$ and all $r \in R$. If $n = 2$, we say f is *R -bilinear*; if $n = 3$, we say f is *R -trilinear*; when R is understood, we may speak of *n -linear* maps.

Let $f : M_1 \times \cdots \times M_n \rightarrow N$ be R -multilinear. By induction on $s \in \mathbb{N}$, we see that

$$f \left(m_1, \dots, m_{i-1}, \sum_{j=1}^s r_j x_j, m_{i+1}, \dots, m_n \right) = \sum_{j=1}^s r_j f(m_1, \dots, x_j, \dots, m_n) \quad (20.3)$$

for all i , where $m_k \in M_k$, $r_j \in R$, and $x_j \in M_i$. Iterating this formula, we obtain

$$\begin{aligned} f & \left(\sum_{j_1=1}^{s_1} r_{1,j_1} x_{1,j_1}, \sum_{j_2=1}^{s_2} r_{2,j_2} x_{2,j_2}, \dots, \sum_{j_n=1}^{s_n} r_{n,j_n} x_{n,j_n} \right) \\ &= \sum_{j_1=1}^{s_1} \sum_{j_2=1}^{s_2} \cdots \sum_{j_n=1}^{s_n} r_{1,j_1} r_{2,j_2} \cdots r_{n,j_n} f(x_{1,j_1}, \dots, x_{n,j_n}). \end{aligned} \quad (20.4)$$

Furthermore, if P is an R -module and $g : N \rightarrow P$ is any R -linear map, then $g \circ f : M_1 \times \cdots \times M_n \rightarrow P$ is R -multilinear. This follows by applying g to each side of (20.1) and (20.2).

20.2 Alternating Maps

For R -modules M and P , we define a map $f : M^n \rightarrow P$ to be *alternating* iff f is R -multilinear and $f(m_1, \dots, m_n) = 0$ whenever $m_i = m_j$ for some $i \neq j$. The alternating condition is related to the following *anti-commutativity conditions* on an R -multilinear map f :

- (AC1) $f(m_1, \dots, m_i, m_{i+1}, \dots, m_n) = -f(m_1, \dots, m_{i+1}, m_i, \dots, m_n)$ for all $i < n$ and all $m_k \in M$. In other words, interchanging two adjacent inputs of f multiplies the output by -1 .
- (AC2) $f(m_1, \dots, m_i, \dots, m_j, \dots, m_n) = -f(m_1, \dots, m_j, \dots, m_i, \dots, m_n)$ for all $i < j$ and all $m_k \in M$. In other words, interchanging any two inputs of f multiplies the output by -1 .
- (AC3) $f(m_{w(1)}, \dots, m_{w(n)}) = \text{sgn}(w)f(m_1, \dots, m_n)$ for all permutations $w \in S_n$ and all $m_k \in M$. In other words, rearranging the inputs of f according to the permutation w multiplies the output by $\text{sgn}(w)$.

We make the following claims regarding these conditions.

Claim 1: The conditions (AC1), (AC2), and (AC3) are equivalent. *Proof:* By letting w be the transposition (i, j) , which satisfies $\text{sgn}((i, j)) = -1$, we see that (AC3) implies (AC2). Evidently (AC2) implies (AC1). To see that (AC1) implies (AC3), recall from Chapter 2 that the list of inputs $(m_{w(1)}, \dots, m_{w(n)})$ can be sorted into the list (m_1, \dots, m_n) using exactly $\text{inv}(w(1), \dots, w(n)) = \text{inv}(w)$ basic transposition moves, where a basic transposition move switches two adjacent elements in a list. According to (AC1), each such move multiplies the value of f by -1 . Therefore,

$$f(m_1, \dots, m_n) = (-1)^{\text{inv}(w)} f(m_{w(1)}, \dots, m_{w(n)}) = \text{sgn}(w)f(m_{w(1)}, \dots, m_{w(n)}).$$

Since $\text{sgn}(w) = \pm 1$, this relation is equivalent to (AC3). We say that an R -multilinear map $f : M^n \rightarrow P$ is *anti-commutative* iff the equivalent conditions (AC1), (AC2), and (AC3) hold for f .

Claim 2: The alternating condition implies all the anti-commutativity conditions. *Proof:* Assume f is alternating; we show that condition (AC1) holds. Fix $i < n$, and fix the arguments of f at positions different from i and $i + 1$. For any $x, y \in M$, the alternating property gives

$$f(\dots, x + y, x + y, \dots) = f(\dots, x, x, \dots) = f(\dots, y, y, \dots) = 0,$$

where the displayed arguments occur at positions i and $i + 1$. On the other hand, linearity of f in its i 'th and $(i + 1)$ 'th arguments shows that

$$\begin{aligned} f(\dots, x + y, x + y, \dots) &= f(\dots, x, x + y, \dots) + f(\dots, y, x + y, \dots) \\ &= f(\dots, x, x, \dots) + f(\dots, x, y, \dots) + f(\dots, y, x, \dots) + f(\dots, y, y, \dots). \end{aligned} \quad (20.5)$$

Substituting zero in three places and rearranging, we get $f(\dots, x, y, \dots) = -f(\dots, y, x, \dots)$, as needed.

Claim 3: For rings R such that $1_R + 1_R$ is not zero and not a zero divisor, any of the anti-commutativity conditions implies the alternating condition. (For instance, the result holds when R is a field or integral domain such that $1_R + 1_R \neq 0_R$.) *Proof:* We deduce the alternating condition from condition (AC2). Suppose $(m_1, \dots, m_n) \in M^n$ is such that $m_i = m_j$ where $i < j$. By (AC2),

$$f(m_1, \dots, m_i, \dots, m_j, \dots, m_n) = -f(m_1, \dots, m_j, \dots, m_i, \dots, m_n).$$

Since $m_i = m_j$, this relation gives

$$f(m_1, \dots, m_i, \dots, m_j, \dots, m_n) = -f(m_1, \dots, m_i, \dots, m_j, \dots, m_n).$$

Grouping terms, $(1_R + 1_R)f(m_1, \dots, m_i, \dots, m_j, \dots, m_n) = 0$ in R . By hypothesis on R , it follows that $f(m_1, \dots, m_n) = 0$.

Claim 4: If $g : P \rightarrow Q$ is R -linear and $f : M^n \rightarrow P$ is alternating (resp. anti-commutative), then $g \circ f : M^n \rightarrow Q$ is alternating (resp. anti-commutative). *Proof:* This follows by applying g to each side of the identities defining the alternating or anti-commutative properties.

20.3 Symmetric Maps

For R -modules M and P , we define a map $f : M^n \rightarrow P$ to be *symmetric* iff f is R -multilinear and $f(m_1, \dots, m_n) = f(m'_1, \dots, m'_n)$ whenever the list $(m_1, \dots, m_n) \in M^n$ is a rearrangement of the list (m'_1, \dots, m'_n) . It is equivalent to require that

$$f(m_1, \dots, m_i, \dots, m_j, \dots, m_n) = f(m_1, \dots, m_j, \dots, m_i, \dots, m_n)$$

for all $i < j$; i.e., the value of f is unchanged whenever two distinct inputs of f are interchanged. It is also equivalent to require that

$$f(m_1, \dots, m_i, m_{i+1}, \dots, m_n) = f(m_1, \dots, m_{i+1}, m_i, \dots, m_n)$$

for all $i < n$; i.e., the value of f is unchanged whenever two adjacent inputs of f are interchanged. The equivalence of the conditions follows from the fact that an arbitrary rearrangement of the list (m_1, \dots, m_n) can be accomplished by a finite sequence of interchanges of two adjacent arguments (see §2.6).

If $g : P \rightarrow Q$ is R -linear and $f : M^n \rightarrow P$ is symmetric, then $g \circ f : M^n \rightarrow Q$ is symmetric. This follows by applying g to each side of the identities $f(m_1, \dots, m_n) = f(m'_1, \dots, m'_n)$ in the definition of a symmetric map.

20.4 Tensor Product of Modules

Suppose P and M_1, M_2, \dots, M_n are R -modules, and let $X = M_1 \times \dots \times M_n$. Recall the distinction between R -linear maps from X to P and R -multilinear maps from X to P (§20.1). It would be convenient if we could somehow reduce the study of multilinear maps to the study of R -linear maps. This suggests the following universal mapping problem.

Problem (UMP for Tensor Products). Given a commutative ring R and R -modules M_1, \dots, M_n , construct an R -module N and an R -multilinear map $j : M_1 \times \dots \times M_n \rightarrow N$ satisfying the following UMP: for any R -module P , there is a bijection from the set

$$A = \{R\text{-linear maps } g : N \rightarrow P\}$$

onto the set

$$B = \{R\text{-multilinear maps } f : M_1 \times \dots \times M_n \rightarrow P\}$$

sending g to $g \circ j$ for all $g \in A$. In other words, for each R -multilinear map $f : M_1 \times \dots \times M_n \rightarrow P$, there exists a unique R -linear map $g : N \rightarrow P$ with $f = g \circ j$:

$$\begin{array}{ccc} M_1 \times \dots \times M_n & \xrightarrow{j} & N \\ & \searrow f & \downarrow g \\ & & P \end{array}$$

Construction of Solution to the UMP. We will define N to be the quotient of a certain free R -module F by a certain submodule K . The main idea of the construction is to fit together the diagrams describing two previously solved universal mapping problems, as shown here and in Figure 20.1:

$$\begin{array}{ccccc} X & \xrightarrow{i} & F & \xrightarrow{\nu} & N \\ & \searrow f & \downarrow h & \swarrow g & \\ & & P & & \end{array}$$

To begin the construction, let X be the set $M_1 \times \dots \times M_n$, let F be a free R -module with basis X (see §17.12), and let $i : X \rightarrow F$ be the inclusion map. Recall that each element of F can be written uniquely as a finite R -linear combination $c_1x_1 + \dots + c_kx_k$, where $c_j \in R$ and $x_j \in X$. Also recall the universal mapping property of i (§19.1): for any R -module P , there is a bijection α from the set

$$C' = \{\text{all } R\text{-linear maps } h : F \rightarrow P\}$$

onto the set

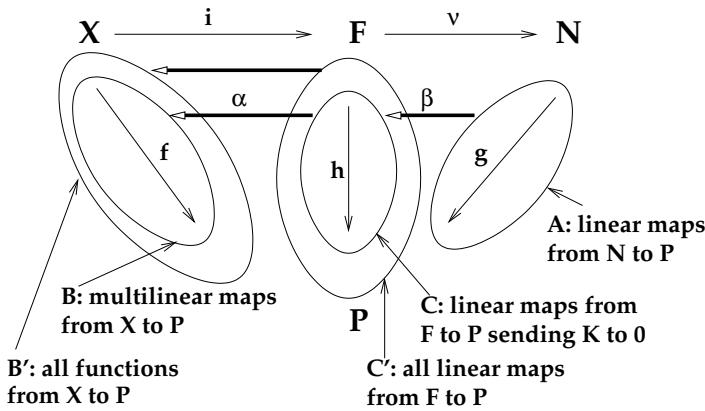
$$B' = \{\text{all functions } f : X \rightarrow P\},$$

given by $\alpha(h) = h \circ i$ for all $h \in C'$.

Next, let K be the R -submodule of F generated by all elements of F of the form

$$\begin{aligned} 1_R(m_1, \dots, m_{k-1}, m_k + m'_k, m_{k+1}, \dots, m_n) \\ - 1_R(m_1, \dots, m_k, \dots, m_n) - 1_R(m_1, \dots, m'_k, \dots, m_n), \end{aligned} \quad (20.6)$$

$$1_R(m_1, \dots, m_{k-1}, rm_k, m_{k+1}, \dots, m_n) - r(m_1, \dots, m_k, \dots, m_n), \quad (20.7)$$

**FIGURE 20.1**

Bijections between Sets of Maps Used to Construct Tensor Products.

where $1 \leq k \leq n$, $m_k, m'_k \in M_k$, $m_s \in M_s$ for $s \neq k$, and $r \in R$. Let N be the R -module F/K , and let $\nu : F \rightarrow N$ be the projection map given by $\nu(z) = z + K$ for $z \in F$. Recall the universal mapping property of ν (§19.2): for any R -module P , there is a bijection β from the set

$$A = \{\text{all } R\text{-linear maps } g : N = F/K \rightarrow P\}$$

onto the set

$$C = \{\text{all } R\text{-linear maps } h : F \rightarrow P \text{ such that } h(z) = 0 \text{ for all } z \in K\},$$

given by $\beta(g) = g \circ \nu$ for all $g \in A$.

Note that $C \subseteq C'$. We claim that $\alpha[C] = B$, the set of all R -multilinear maps from $X = M_1 \times \dots \times M_n$ to P . *Proof:* An R -linear map $h : F \rightarrow P$ in C' belongs to C iff $h[K] = \{0\}$ iff h maps every generator of the submodule K to zero iff

$$\begin{aligned} h(m_1, \dots, m_{k-1}, m_k + m'_k, m_{k+1}, \dots, m_n) \\ - h(m_1, \dots, m_k, \dots, m_n) - h(m_1, \dots, m'_k, \dots, m_n) = 0, \end{aligned}$$

$$h(m_1, \dots, m_{k-1}, rm_k, m_{k+1}, \dots, m_n) - rh(m_1, \dots, m_k, \dots, m_n) = 0$$

for all choices of the variables iff

$$\begin{aligned} (h \circ i)(m_1, \dots, m_k + m'_k, \dots, m_n) &= (h \circ i)(m_1, \dots, m_k, \dots, m_n) + (h \circ i)(m_1, \dots, m'_k, \dots, m_n), \\ (h \circ i)(m_1, \dots, rm_k, \dots, m_n) &= r(h \circ i)(m_1, \dots, m_k, \dots, m_n) \end{aligned}$$

for all choices of the variables iff $\alpha(h) = h \circ i : X \rightarrow P$ is R -multilinear iff $\alpha(h) \in B$. By the claim, the restriction $\alpha|_C : C \rightarrow B$ is a bijection sending h to $h \circ i$ for $h \in C$. We also have the bijection $\beta : A \rightarrow C$ given by $\beta(g) = g \circ \nu$ for $g \in A$. Composing these bijections, we obtain a bijection γ from A to B given by $\gamma(g) = g \circ (\nu \circ i)$ for $g \in A$. Letting $j = \nu \circ i : X \rightarrow N$, we see that $\gamma : A \rightarrow B$ is given by composition with j . See Figure 20.1.

To summarize the construction, we have $N = F/K$, where F is the free R -module with basis $X = M_1 \times \dots \times M_n$ and K is the R -submodule generated by all elements of the form (20.6) and (20.7). The map $j : X \rightarrow N$ sends $(m_1, \dots, m_n) \in X$ to the coset $(m_1, \dots, m_n) + K$ in N . By the coset equality theorem, we know that

$$(m_1, \dots, m_k + m'_k, \dots, m_n) + K = [(m_1, \dots, m_k, \dots, m_n) + K] + [(m_1, \dots, m'_k, \dots, m_n) + K],$$

$$(m_1, \dots, rm_k, \dots, m_n) + K = r[(m_1, \dots, m_k, \dots, m_n) + K],$$

and these observations show that $j : X \rightarrow N$ is an R -multilinear map from X to N , as needed. We call N the *tensor product over R of the modules M_1, \dots, M_n* and write

$$N = M_1 \otimes_R M_2 \otimes_R \cdots \otimes_R M_n.$$

Given $m_k \in M_k$, we introduce the *tensor notation* $m_1 \otimes m_2 \otimes \cdots \otimes m_n = j(m_1, \dots, m_n) \in N$. In this notation, R -multilinearity of j translates into the identities

$$m_1 \otimes \cdots \otimes (m_k + m'_k) \otimes \cdots \otimes m_n = m_1 \otimes \cdots \otimes m_k \otimes \cdots \otimes m_n + m_1 \otimes \cdots \otimes m'_k \otimes \cdots \otimes m_n; \quad (20.8)$$

$$m_1 \otimes \cdots \otimes (rm_k) \otimes \cdots \otimes m_n = r(m_1 \otimes \cdots \otimes m_k \otimes \cdots \otimes m_n), \quad (20.9)$$

valid for all $m_k, m'_k \in M_k$, $m_s \in M_s$, and $r \in R$.

Uniqueness of Solution to the UMP. To justify our new notation for N and j , we show that the solution (N, j) to our universal mapping problem is unique up to a unique isomorphism compatible with the universal map j . Suppose (N', j') is another solution to the UMP. The proof involves the following four diagrams of sets and mappings:

$$\begin{array}{cccc} \begin{array}{c} X \xrightarrow{j} N \\ \searrow j' \quad \downarrow g \\ N' \end{array} & \begin{array}{c} X \xrightarrow{j'} N' \\ \searrow j \quad \downarrow g' \\ N \end{array} & \begin{array}{c} X \xrightarrow{j} N \\ \searrow j \quad \downarrow h \\ N \end{array} & \begin{array}{c} X \xrightarrow{j'} N' \\ \searrow j' \quad \downarrow h' \\ N' \end{array} \end{array}$$

Since $j' : X \rightarrow N'$ is R -multilinear and (N, j) solves the UMP, we get a unique R -linear map $g : N \rightarrow N'$ with $j' = g \circ j$ (see the first diagram above). It now suffices to show that g is an isomorphism. Since $j : X \rightarrow N$ is R -multilinear and (N', j') solves the UMP, we get a unique R -linear map $g' : N' \rightarrow N$ with $j = g' \circ j'$ (see the second diagram above). It follows that $\text{id}_N \circ j = j = (g' \circ g) \circ j$. By the uniqueness assertion in the UMP for (N, j) , there is only one R -linear map $h : N \rightarrow N$ with $j = h \circ j$ (see the third diagram above). Therefore, $g' \circ g = \text{id}_N$. Similarly, using the uniqueness of h' in the fourth diagram, we see that $g \circ g' = \text{id}_{N'}$. So g' is the two-sided inverse of g , hence both maps are isomorphisms. Note that this uniqueness proof is essentially identical to the earlier uniqueness proof given for the coproduct of a family of sets (§19.7). In general, this same proof template can be used over and over again to establish the uniqueness of solutions of universal mapping problems, up to a unique isomorphism compatible with the universal maps. For future universal mapping problems, we will omit the details of this uniqueness proof, allowing the reader to verify the applicability of the proof template used here.

Now let M be an R -module and n a positive integer. Consider the Cartesian product $M^n = M \times M \times \cdots \times M$, where there are n copies of M . By letting each $M_k = M$ in the tensor product construction, we obtain the n^{th} *tensor power* of M , denoted

$$\bigotimes^n M = M^{\otimes n} = M \otimes_R M \otimes_R \cdots \otimes_R M,$$

and the universal map j that sends $(m_1, \dots, m_n) \in M^n$ to $m_1 \otimes \cdots \otimes m_n \in M^{\otimes n}$. The UMP for tensor products says that there is a bijection γ from the set of R -linear maps (R -module homomorphisms) $g : M^{\otimes n} \rightarrow P$ onto the set of R -multilinear maps $f : M^n \rightarrow P$, given by $\gamma(g) = g \circ j$. In the next two sections, we use the tensor power $M^{\otimes n}$ to solve universal mapping problems for alternating maps and symmetric maps.

20.5 Exterior Powers of a Module

In the last section, we solved a universal mapping problem that converted *R-multilinear* maps into *R-linear* maps. Here, we pose and solve a similar problem that will convert *alternating* maps into *R-linear* maps.

Problem (UMP for Exterior Powers). Given a commutative ring R , an R -module M , and a positive integer n , construct an R -module N and an alternating map $i : M^n \rightarrow N$ satisfying the following UMP: for any R -module P , there is a bijection from the set

$$A = \{R\text{-linear maps } g : N \rightarrow P\}$$

onto the set

$$B = \{\text{alternating maps } f : M^n \rightarrow P\}$$

that sends $g \in A$ to $g \circ i \in B$. In other words, for each alternating map $f : M^n \rightarrow P$, there exists a unique R -linear map $g : N \rightarrow P$ with $f = g \circ i$:

$$\begin{array}{ccc} M^n & \xrightarrow{i} & N \\ & \searrow f & \downarrow g \\ & & P \end{array}$$

Construction of Solution to the UMP. As before, the idea is to fit together two previously solved universal mapping problems, as suggested in the following diagram and in Figure 20.2:

$$\begin{array}{ccccc} M^n & \xrightarrow{j} & M^{\otimes n} & \xrightarrow{\nu} & N \\ & \searrow f & \downarrow h & \swarrow g & \\ & & P & & \end{array}$$

To explain this, first recall the universal mapping property for $M^{\otimes n}$ and $j : M^n \rightarrow M^{\otimes n}$ (§20.4): for any R -module P , there is a bijection α from the set

$$C' = \{\text{all } R\text{-linear maps } h : M^{\otimes n} \rightarrow P\}$$

onto the set

$$B' = \{\text{all } R\text{-multilinear maps } f : M^n \rightarrow P\},$$

given by $\alpha(h) = h \circ j$ for all $h \in C'$.

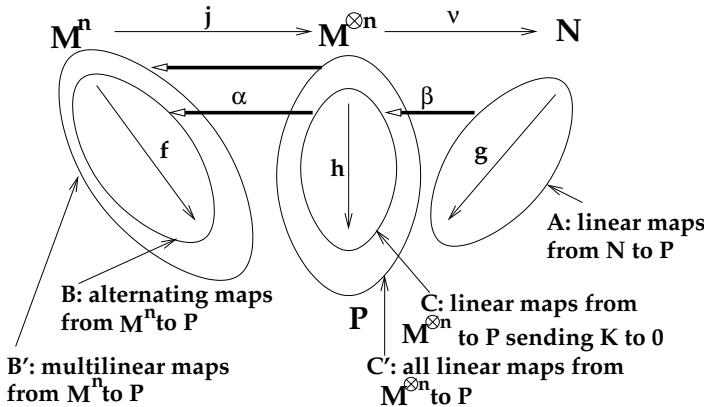
Next, let K be the R -submodule of $M^{\otimes n}$ generated by all elements $m_1 \otimes m_2 \otimes \cdots \otimes m_n \in M^{\otimes n}$ such that $m_k = m_\ell$ for some $k \neq \ell$. Let N be the R -module $M^{\otimes n}/K$, and let $\nu : M^{\otimes n} \rightarrow N$ be the projection map given by $\nu(z) = z + K$ for $z \in M^{\otimes n}$. Recall the universal mapping property of ν (§19.2): for any R -module P , there is a bijection β from the set

$$A = \{\text{all } R\text{-linear maps } g : N = M^{\otimes n}/K \rightarrow P\}$$

onto the set

$$C = \{\text{all } R\text{-linear maps } h : M^{\otimes n} \rightarrow P \text{ such that } h(z) = 0 \text{ for all } z \in K\},$$

given by $\beta(g) = g \circ \nu$ for all $g \in A$.

**FIGURE 20.2**

Bijections between Sets of Maps Used to Construct Exterior Powers.

Note that $C \subseteq C'$. We claim that $\alpha[C] = B$, the set of all alternating maps from M^n to P . *Proof:* An R -linear map $h : M^{\otimes n} \rightarrow P$ in C' belongs to C iff $h[K] = \{0\}$ iff h maps every generator of the submodule K to zero iff $h(m_1 \otimes \cdots \otimes m_n) = 0$ whenever $m_k = m_\ell$ for some $k \neq \ell$ iff $(h \circ j)(m_1, \dots, m_n) = 0$ whenever $m_k = m_\ell$ for some $k \neq \ell$ iff $\alpha(h) = h \circ j : M^n \rightarrow P$ is alternating iff $\alpha(h) \in B$. By the claim, α restricts to a bijection $\alpha|_C : C \rightarrow B$ sending $h \in C$ to $h \circ j \in B$. We also have the bijection $\beta : A \rightarrow C$ sending $g \in A$ to $g \circ \nu \in C$. Composing these bijections, we obtain a bijection $\gamma : A \rightarrow B$ given by $\gamma(g) = g \circ (\nu \circ j)$ for $g \in A$. Letting $i = \nu \circ j : M^n \rightarrow N$, we see that the bijection from A to B is given by composition with i . Since j is R -multilinear and ν is R -linear, the composite map i is R -multilinear. See Figure 20.2.

The map $i : M^n \rightarrow N$ sends $(m_1, \dots, m_n) \in M^n$ to the coset $(m_1 \otimes \cdots \otimes m_n) + K$ in N . By definition of K , $i(m_1, \dots, m_n) = 0 + K = 0_N$ whenever $m_k = m_\ell$ for some $k \neq \ell$. Therefore, i is alternating. The standard argument proves that the solution (N, i) to the UMP is unique up to a unique R -isomorphism.

We call N the n 'th exterior power of M and write $N = \bigwedge^n M$. We also write

$$m_1 \wedge m_2 \wedge \cdots \wedge m_n = i(m_1, \dots, m_n) = m_1 \otimes \cdots \otimes m_n + K \in N$$

and call this element the *wedge product* of m_1, \dots, m_n . In this notation, the R -multilinearity and alternating properties of i translate into the identities:

$$m_1 \wedge \cdots \wedge (m_k + m'_k) \wedge \cdots \wedge m_n = m_1 \wedge \cdots \wedge m_k \wedge \cdots \wedge m_n + m_1 \wedge \cdots \wedge m'_k \wedge \cdots \wedge m_n;$$

$$m_1 \wedge \cdots \wedge (rm_k) \wedge \cdots \wedge m_n = r(m_1 \wedge \cdots \wedge m_k \wedge \cdots \wedge m_n);$$

$$m_1 \wedge \cdots \wedge m_n = 0 \text{ whenever } m_k = m_\ell \text{ for some } k \neq \ell.$$

The anti-commutativity of i , which follows from the alternating property, translates into the following facts:

$$m_1 \wedge \cdots \wedge m_k \wedge m_{k+1} \wedge \cdots \wedge m_n = -m_1 \wedge \cdots \wedge m_{k+1} \wedge m_k \wedge \cdots \wedge m_n;$$

$$m_1 \wedge \cdots \wedge m_k \wedge \cdots \wedge m_\ell \wedge \cdots \wedge m_n = -m_1 \wedge \cdots \wedge m_\ell \wedge \cdots \wedge m_k \wedge \cdots \wedge m_n;$$

$$m_{w(1)} \wedge \cdots \wedge m_{w(n)} = (\operatorname{sgn}(w))m_1 \wedge \cdots \wedge m_n \text{ for all } w \in S_n.$$

20.6 Symmetric Powers of a Module

Next we devise a universal construction for converting *symmetric* maps into *R-linear* maps.

Problem (UMP for Symmetric Powers). Given a commutative ring R , an R -module M , and a positive integer n , construct an R -module N and a symmetric map $i : M^n \rightarrow N$ satisfying the following UMP: for any R -module P , there is a bijection from the set

$$A = \{R\text{-linear maps } g : N \rightarrow P\}$$

onto the set

$$B = \{\text{symmetric maps } f : M^n \rightarrow P\}$$

sending $g \in A$ to $g \circ i \in B$. In other words, for each symmetric map $f : M^n \rightarrow P$, there exists a unique R -linear map $g : N \rightarrow P$ with $f = g \circ i$:

$$\begin{array}{ccc} M^n & \xrightarrow{i} & N \\ & \searrow f & \downarrow g \\ & & P \end{array}$$

Construction of Solution to the UMP. The proof is nearly identical to what we did for alternating maps (see the diagram below and Figure 20.3).

$$\begin{array}{ccccc} M^n & \xrightarrow{j} & M^{\otimes n} & \xrightarrow{\nu} & N \\ & \searrow f & \downarrow h & \nearrow g & \\ & & P & & \end{array}$$

To start, recall once again the universal mapping property for $M^{\otimes n}$ and $j : M^n \rightarrow M^{\otimes n}$ (§20.4): for any R -module P , there is a bijection α from the set

$$C' = \{\text{all } R\text{-linear maps } h : M^{\otimes n} \rightarrow P\}$$

onto the set

$$B' = \{\text{all } R\text{-multilinear maps } f : M^n \rightarrow P\},$$

given by $\alpha(h) = h \circ j$ for all $h \in C'$.

Next, let K be the R -submodule of $M^{\otimes n}$ generated by all elements of the form

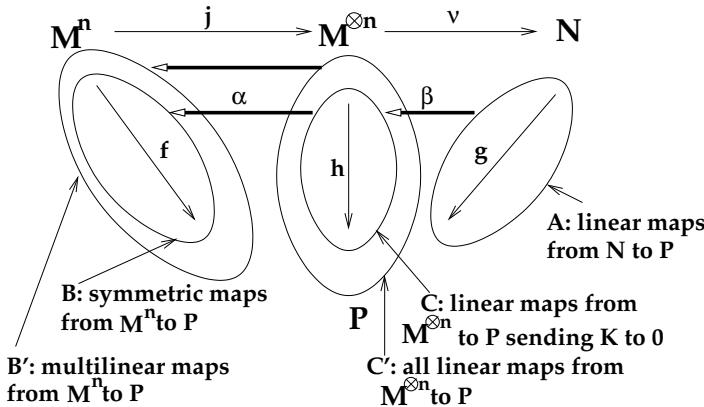
$$m_1 \otimes m_2 \otimes \cdots \otimes m_n - m'_1 \otimes m'_2 \otimes \cdots \otimes m'_n,$$

where the list (m'_1, \dots, m'_n) is a rearrangement of the list (m_1, \dots, m_n) . Let N be the R -module $M^{\otimes n}/K$, and let $\nu : M^{\otimes n} \rightarrow N$ be the projection map given by $\nu(z) = z + K$ for $z \in M^{\otimes n}$. Recall once again the universal mapping property of ν (§19.2): for any R -module P , there is a bijection β from the set

$$A = \{\text{all } R\text{-linear maps } g : N = M^{\otimes n}/K \rightarrow P\}$$

onto the set

$$C = \{\text{all } R\text{-linear maps } h : M^{\otimes n} \rightarrow P \text{ such that } h(z) = 0 \text{ for all } z \in K\},$$

**FIGURE 20.3**

Bijections between Sets of Maps Used to Construct Symmetric Powers.

given by $\beta(g) = g \circ \nu$ for all $g \in A$.

Note that $C \subseteq C'$. We claim that $\alpha[C] = B$, the set of all symmetric maps from M^n to P . *Proof:* An R -linear map $h : M^{\otimes n} \rightarrow P$ in C' belongs to C iff $h[K] = \{0\}$ iff h maps every generator of the submodule K to zero iff $h(m_1 \otimes \cdots \otimes m_n) - h(m'_1 \otimes \cdots \otimes m'_n) = 0$ whenever (m'_1, \dots, m'_n) is a rearrangement of (m_1, \dots, m_n) iff $(h \circ j)(m_1, \dots, m_n) = (h \circ j)(m'_1, \dots, m'_n)$ whenever (m'_1, \dots, m'_n) is a rearrangement of (m_1, \dots, m_n) iff $\alpha(h) = h \circ j : M^n \rightarrow P$ is symmetric iff $\alpha(h) \in B$. By the claim, α restricts to a bijection $\alpha|_C : C \rightarrow B$ sending $h \in C$ to $h \circ j \in B$. We also have the bijection $\beta : A \rightarrow C$ sending $g \in A$ to $g \circ \nu \in C$. Composing these bijections, we obtain a bijection $\gamma : A \rightarrow B$ given by $\gamma(g) = g \circ (\nu \circ j)$ for $g \in A$. Letting $i = \nu \circ j : M^n \rightarrow N$, we see that the bijection from A to B is given by composition with i . Since j is R -multilinear and ν is R -linear, the composite map i is R -multilinear. See Figure 20.3.

The map $i : M^n \rightarrow N$ sends $(m_1, \dots, m_n) \in M^n$ to the coset $(m_1 \otimes \cdots \otimes m_n) + K$ in N . By definition of K , $i(m_1, \dots, m_n) = i(m'_1, \dots, m'_n)$ in N whenever (m'_1, \dots, m'_n) is a rearrangement of (m_1, \dots, m_n) . Therefore, i is symmetric. The standard argument proves that the solution (N, i) to the UMP is unique up to a unique R -isomorphism.

We call N the n 'th symmetric power of M and write $N = \text{Sym}^n M$. We also write

$$m_1 m_2 \cdots m_n = i(m_1, \dots, m_n) = m_1 \otimes \cdots \otimes m_n + K \in N$$

and call this element the *symmetric product* of m_1, \dots, m_n . In this notation, the R -multilinearity and symmetric properties of i translate into the identities:

$$m_1 m_2 \cdots (m_k + m'_k) \cdots m_n = m_1 m_2 \cdots m_k \cdots m_n + m_1 m_2 \cdots m'_k \cdots m_n;$$

$$m_1 m_2 \cdots (r m_k) \cdots m_n = r(m_1 m_2 \cdots m_k \cdots m_n);$$

$$m_1 m_2 \cdots m_n = m'_1 m'_2 \cdots m'_n$$

whenever (m'_1, \dots, m'_n) is a rearrangement of (m_1, \dots, m_n) .

20.7 Myths about Tensor Products

The tensor product construction is more subtle than other module constructions such as the direct product. To help the reader avoid common errors, we discuss some misconceptions or “myths” about tensor products in this section. For simplicity, let us consider the tensor product $M \otimes_R N$ of two R -modules M and N .

Myth 1: “Every element of $M \otimes_R N$ has the form $x \otimes y$ for some $x \in M$ and $y \in N$.” This myth arises from the faulty assumption that the underlying set of the module $M \otimes_R N$ is the Cartesian product $M \times N = \{(x, y) : x \in M, y \in N\}$. However, recalling the construction in §20.4, we see that $M \otimes_R N$ was built by starting with a free module F having the set $M \times N$ as a basis, and then taking the quotient by a certain submodule K . For any basis element (x, y) of F , we wrote $x \otimes y$ for the coset $(x, y) + K$. Since every element of F is a finite R -linear combination of basis elements, it is true that every element of $M \otimes_R N$ is a finite R -linear combination of the form $\sum_{i=1}^k r_i(x_i \otimes y_i)$, where $k \in \mathbb{N}$, $r_i \in R$, $x_i \in M$, and $y_i \in N$. In fact, we can use (20.9) to write $r_i(x_i \otimes y_i) = (r_i x_i) \otimes y_i$ for each i . This means that *every element of $M \otimes_R N$ can be written as a finite sum of “basic” tensors $x \otimes y$ with $x \in M$ and $y \in N$.*

Myth 2: “The set $\{x \otimes y : x \in M, y \in N\}$ is a basis for the R -module $M \otimes_R N$.” In fact, the set of all ordered pairs (x, y) with $x \in M$ and $y \in N$ is a basis for the free R -module F used in the construction of $M \otimes_R N$. However, after passing to the quotient module F/K , the set of cosets $(x, y) + K$ is almost never a basis for F/K . For instance, using (20.8), we have $(x_1 + x_2) \otimes y - (x_1 \otimes y) - (x_2 \otimes y) = 0$ for all $x_1, x_2 \in M$ and $y \in N$, which gives a dependence relation among three basic tensors. On the other hand, as we saw above, it is true that the set of all tensors $x \otimes y$ forms a *generating set* for the R -module $M \otimes_R N$. In the coming sections, we will find bases for tensor products, exterior powers, and symmetric powers of free modules.

Myth 3: “All we need to do to define a function with domain $M \otimes_R N$ is to declare $f(x \otimes y)$ to be any formula involving x and y .” To illustrate the problems that can arise here, consider the following proposed proof that the R -modules $M \otimes_R N$ and $N \otimes_R M$ are isomorphic: “Define $f : M \otimes_R N \rightarrow N \otimes_R M$ by $f(x \otimes y) = y \otimes x$ for all $x \in M$ and $y \in N$. The map f is evidently an R -linear bijection, so $M \otimes_R N \cong N \otimes_R M$.” The first difficulty is that the stated formula does not define f on the entire domain $M \otimes_R N$ (see Myth 1). Since most maps of interest will preserve addition, one could try to get around this by “extending f additively,” i.e., by setting $f(\sum_{i=1}^k x_i \otimes y_i) = \sum_{i=1}^k y_i \otimes x_i$.

However, there is still a problem with this extended definition. When computing $f(z)$ for any $z \in M \otimes_R N$, the output appears to depend on the particular representation of z as a sum of tensors $x_i \otimes y_i$. Usually, z has many such representations; how do we know that the formula for f will give the same answer no matter which representation we use? For similar reasons, it should now be far from obvious that f must be one-to-one.

To resolve these difficulties in defining f , we must return to the universal mapping property characterizing $M \otimes_R N$. We have seen that there is a bijection between R -linear maps from $M \otimes_R N$ to any other fixed module P and R -bilinear maps from $M \times N$ to P . In this case, P is the R -module $N \otimes_R M$. To define the required R -linear map $f : M \otimes_R N \rightarrow N \otimes_R M$, we instead define a function $g : M \times N \rightarrow N \otimes_R M$ on the product set $M \times N$, by letting $g(x, y) = y \otimes x$ for $x \in M$ and $y \in N$. Note that this definition has none of the problems that we encountered earlier, since every element of the Cartesian product $M \times N$ can be written uniquely as (x, y) for some $x \in M$ and $y \in N$. Furthermore, using (20.8) and (20.9) in $N \otimes_R M$, we see that g is R -bilinear: for all $x, x_1, x_2 \in M$ and

$y, y_1, y_2 \in N$ and $r \in R$,

$$g(x_1 + x_2, y) = y \otimes (x_1 + x_2) = (y \otimes x_1) + (y \otimes x_2) = g(x_1, y) + g(x_2, y);$$

$$g(x, y_1 + y_2) = (y_1 + y_2) \otimes x = (y_1 \otimes x) + (y_2 \otimes x) = g(x, y_1) + g(x, y_2);$$

$$g(rx, y) = y \otimes (rx) = r(y \otimes x) = rg(x, y); \quad g(x, ry) = (ry) \otimes x = r(y \otimes x) = rg(x, y).$$

The UMP now provides us with a unique (and well-defined) R -linear map $f : M \otimes_R N \rightarrow N \otimes_R M$ with $g = f \circ j$; i.e., $f(x \otimes y) = f(j(x, y)) = g(x, y) = y \otimes x$. This is our original “definition” of f on generators of $M \otimes_R N$, but now we are sure that f is well-defined and R -linear. In most situations, this is the method one must use to define maps whose domain is a tensor product, exterior power, or symmetric power.

Now, why is f a bijection? Seeing that f is onto is not too hard, but checking the injectivity of f directly can be difficult. Instead, we show that f is bijective by exhibiting a two-sided inverse map. Start with the map $g_1 : N \times M \rightarrow M \otimes_R N$ given by $g_1(y, x) = x \otimes y$ for all $y \in N$ and $x \in M$. As above, we see that g_1 is R -bilinear, so the UMP furnishes a unique R -linear map $f_1 : N \otimes_R M \rightarrow M \otimes_R N$ given on generators by $f_1(y \otimes x) = x \otimes y$. To check that $f \circ f_1 = \text{id}_{N \otimes_R M}$, observe that $f \circ f_1(y \otimes x) = f(x \otimes y) = y \otimes x$ for every generator $y \otimes x$ of the R -module $N \otimes_R M$. The R -linear map $\text{id}_{N \otimes_R M}$ also sends $y \otimes x$ to $y \otimes x$ for all $y \in N$ and $x \in M$. We know that two R -linear maps agreeing on a set of generators for their common domain must be equal. So $f \circ f_1 = \text{id}_{N \otimes_R M}$, and similarly $f_1 \circ f = \text{id}_{M \otimes_R N}$. Finally, we have a rigorous proof of the module isomorphism $M \otimes_R N \cong N \otimes_R M$.

20.8 Tensor Product Isomorphisms

In this section, we give further illustrations of the technique used above to prove that $M \otimes_R N$ is isomorphic to $N \otimes_R M$. As a first example, we prove that *for every R -module M , $R \otimes_R M \cong M$* .

Given an R -module M , let $r \star m$ denote the action of a scalar $r \in R$ on a module element $m \in M$. Define a map $g : R \times M \rightarrow M$ by $g(r, m) = r \star m$ for $r \in R$ and $m \in M$. Now, g is R -bilinear, since for $r, s \in R$ and $m, n \in M$:

$$\begin{aligned} g(r+s, m) &= (r+s) \star m = (r \star m) + (s \star m) = g(r, m) + g(s, m); \\ g(rs, m) &= (rs) \star m = r \star (s \star m) = rg(s, m); \\ g(r, m+n) &= r \star (m+n) = (r \star m) + (r \star n) = g(r, m) + g(r, n); \\ g(r, sm) &= r \star (sm) = (rs) \star m = (sr) \star m = s \star (r \star m) = sg(r, m). \end{aligned}$$

(Note that commutativity of R was needed to justify the last equality.) We therefore get a unique R -linear map $f : R \otimes_R M \rightarrow M$ defined on generators by $f(r \otimes m) = r \star m$ for $r \in R$ and $m \in M$.

Next we define $h : M \rightarrow R \otimes_R M$ by $h(m) = 1_R \otimes m$ for $m \in M$. Now, h is R -linear because $h(m+n) = 1 \otimes (m+n) = (1 \otimes m) + (1 \otimes n) = h(m) + h(n)$ for all $m, n \in M$; and $h(rm) = 1 \otimes (rm) = r(1 \otimes m) = rh(m)$ for all $r \in R$ and $m \in M$. We claim h is the two-sided inverse of f , so that both maps are R -module isomorphisms. On one hand, for any $m \in M$, $f(h(m)) = f(1 \otimes m) = 1 \star m = m = \text{id}_M(m)$, so $f \circ h = \text{id}_M$. On the other hand, for any $r \in R$ and $m \in M$,

$$h(f(r \otimes m)) = h(r \star m) = 1 \otimes (rm) = r(1 \otimes m) = (r1) \otimes m = r \otimes m.$$

Thus, the R -linear maps $h \circ f$ and $\text{id}_{R \otimes_R M}$ have the same effect on all generators of $R \otimes_R M$, so these maps are equal. This completes the proof that $R \otimes_R M \cong M$ as R -modules.

For our next result, fix a ring R and R -modules M , N , and P . We will prove that

$$(M \oplus N) \otimes_R P \cong (M \otimes_R P) \oplus (N \otimes_R P)$$

as R -modules. (Recall that $M \oplus N$ is the set of ordered pairs (m, n) with $m \in M$ and $n \in N$, which becomes an R -module under componentwise operations.)

First, define a map $g : (M \oplus N) \times P \rightarrow (M \otimes_R P) \oplus (N \otimes_R P)$ by $g((m, n), p) = (m \otimes p, n \otimes p)$ for all $m \in M$, $n \in N$, and $p \in P$. One checks that g is R -bilinear; for example, if $u = (m_1, n_1)$ and $v = (m_2, n_2)$ are in $M \oplus N$, we compute:

$$\begin{aligned} g(u + v, p) &= g((m_1 + m_2, n_1 + n_2), p) = ((m_1 + m_2) \otimes p, (n_1 + n_2) \otimes p) \\ &= (m_1 \otimes p + m_2 \otimes p, n_1 \otimes p + n_2 \otimes p) \\ &= (m_1 \otimes p, n_1 \otimes p) + (m_2 \otimes p, n_2 \otimes p) \\ &= g((m_1, n_1), p) + g((m_2, n_2), p) = g(u, p) + g(v, p). \end{aligned}$$

By the UMP, there is a unique R -linear map $f : (M \oplus N) \otimes_R P \rightarrow (M \otimes_R P) \oplus (N \otimes_R P)$ given on generators by

$$f((m, n) \otimes p) = (m \otimes p, n \otimes p).$$

To construct the inverse of f , we first define maps $g_1 : M \times P \rightarrow (M \oplus N) \otimes_R P$ and $g_2 : N \times P \rightarrow (M \oplus N) \otimes_R P$, by setting $g_1(m, p) = (m, 0) \otimes p$ and $g_2(n, p) = (0, n) \otimes p$ for all $m \in M$, $n \in N$, and $p \in P$. One quickly verifies that g_1 and g_2 are R -bilinear, so the UMP provides R -linear maps $h_1 : M \otimes_R P \rightarrow (M \oplus N) \otimes_R P$ and $h_2 : N \otimes_R P \rightarrow (M \oplus N) \otimes_R P$ defined on generators by

$$h_1(m \otimes p) = (m, 0) \otimes p, \quad h_2(n \otimes p) = (0, n) \otimes p.$$

Now, using the UMP for the direct sum of two R -modules (§19.4), we can combine h_1 and h_2 to get an R -linear map $h : (M \otimes_R P) \oplus (N \otimes_R P) \rightarrow (M \oplus N) \otimes_R P$ given by $h(y, z) = h_1(y) + h_2(z)$ for $y \in M \otimes_R P$ and $z \in N \otimes_R P$.

To finish the proof, we prove that $f \circ h$ and $h \circ f$ are identity maps by checking that each map sends all relevant generators to themselves. A typical generator of $(M \oplus N) \otimes_R P$ is $w = (m, n) \otimes p$ with $m \in M$, $n \in N$, and $p \in P$. We compute

$$\begin{aligned} h(f(w)) &= h(m \otimes p, n \otimes p) = h_1(m \otimes p) + h_2(n \otimes p) \\ &= (m, 0) \otimes p + (0, n) \otimes p = ((m, 0) + (0, n)) \otimes p = (m, n) \otimes p = w. \end{aligned}$$

On the other hand, one readily checks that elements of the form $(m \otimes p, 0)$ and $(0, n \otimes p)$ (with $m \in M$, $n \in N$, and $p \in P$) generate the R -module $(M \otimes_R P) \oplus (N \otimes_R P)$. For generators of the first type,

$$f(h(m \otimes p, 0)) = f(h_1(m \otimes p)) = f((m, 0) \otimes p) = (m \otimes p, 0 \otimes p) = (m \otimes p, 0).$$

For generators of the second type,

$$f(h(0, n \otimes p)) = f(h_2(n \otimes p)) = f((0, n) \otimes p) = (0 \otimes p, n \otimes p) = (0, n \otimes p).$$

(These calculations use the readily verified fact that $0 \otimes p = 0$ for any $p \in P$.) We conclude that h is the inverse of f , so both maps are R -module isomorphisms.

An entirely analogous proof shows that $P \otimes_R (M \oplus N) \cong (P \otimes_R M) \oplus (P \otimes_R N)$. More generally, for all R -modules M_i and P_j , the same techniques (Exercise 31) prove that

$$\left(\bigoplus_{i \in I} M_i \right) \otimes_R P_2 \otimes_R \cdots \otimes_R P_k \cong \bigoplus_{i \in I} (M_i \otimes_R P_2 \otimes_R \cdots \otimes_R P_k). \quad (20.10)$$

Even more generally, there is an isomorphism

$$\left(\bigoplus_{i_1 \in I_1} M_{i_1} \right) \otimes_R \cdots \otimes_R \left(\bigoplus_{i_k \in I_k} M_{i_k} \right) \cong \bigoplus_{(i_1, \dots, i_k) \in I_1 \times \cdots \times I_k} (M_{i_1} \otimes_R M_{i_2} \otimes_R \cdots \otimes_R M_{i_k}). \quad (20.11)$$

20.9 Associativity of Tensor Products

Next we prove an associativity result for tensor products. Fix an R -module M and positive integers k and m . We will show that $M^{\otimes k} \otimes_R M^{\otimes m} \cong M^{\otimes(k+m)}$ as R -modules by constructing isomorphisms in both directions. First, define $g : M^{k+m} \rightarrow M^{\otimes k} \otimes_R M^{\otimes m}$ by

$$g(x_1, \dots, x_k, x_{k+1}, \dots, x_{k+m}) = (x_1 \otimes \cdots \otimes x_k) \otimes (x_{k+1} \otimes \cdots \otimes x_{k+m})$$

for $x_i \in M$. It is routine to check (by repeated use of (20.8) and (20.9)) that g is an R -multilinear map. Hence g induces a unique R -linear map $f : M^{\otimes(k+m)} \rightarrow M^{\otimes k} \otimes_R M^{\otimes m}$ defined on generators by

$$f(x_1 \otimes \cdots \otimes x_k \otimes x_{k+1} \otimes \cdots \otimes x_{k+m}) = (x_1 \otimes \cdots \otimes x_k) \otimes (x_{k+1} \otimes \cdots \otimes x_{k+m}).$$

Constructing f^{-1} is somewhat tricky, since the domain of f^{-1} involves three tensor products. To begin, fix $z = (x_{k+1}, \dots, x_{k+m}) \in M^m$ and define a map $p'_z : M^k \rightarrow M^{\otimes(k+m)}$ by

$$p'_z(x_1, \dots, x_k) = x_1 \otimes \cdots \otimes x_k \otimes x_{k+1} \otimes \cdots \otimes x_{k+m}$$

for all $(x_1, \dots, x_k) \in M^k$. One checks that p'_z is k -linear. Invoking the UMP for $M^{\otimes k}$, we get an R -linear map $p_z : M^{\otimes k} \rightarrow M^{\otimes(k+m)}$ given on generators by

$$p_z(x_1 \otimes \cdots \otimes x_k) = x_1 \otimes \cdots \otimes x_k \otimes x_{k+1} \otimes \cdots \otimes x_{k+m}.$$

Next, for each fixed $y \in M^{\otimes k}$, we define a map $q'_y : M^m \rightarrow M^{\otimes(k+m)}$ by $q'_y(z) = p_z(y)$ for all $z \in M^m$. One checks, using the formula for p_z on generators, that q'_y is m -linear. Invoking the UMP for $M^{\otimes m}$, we get an R -linear map $q_y : M^{\otimes m} \rightarrow M^{\otimes(k+m)}$ given on generators by

$$q_y(x_{k+1} \otimes \cdots \otimes x_{k+m}) = p_{(x_{k+1}, \dots, x_{k+m})}(y).$$

Finally, define a map $t' : M^{\otimes k} \times M^{\otimes m} \rightarrow M^{\otimes(k+m)}$ by setting $t'(y, w) = q_y(w)$ for all $y \in M^{\otimes k}$ and all $w \in M^{\otimes m}$. One checks that t' is R -bilinear, so we finally get an R -linear map $t : M^{\otimes k} \otimes_R M^{\otimes m} \rightarrow M^{\otimes(k+m)}$ given on generators by $t(y \otimes w) = q_y(w)$. Tracing through all the definitions, we find that

$$t((x_1 \otimes \cdots \otimes x_k) \otimes (x_{k+1} \otimes \cdots \otimes x_{k+m})) = x_1 \otimes \cdots \otimes x_k \otimes x_{k+1} \otimes \cdots \otimes x_{k+m}$$

for all $x_j \in M$. Hence, $t \circ f = \text{id}_{M^{\otimes(k+m)}}$ since these two R -linear maps agree on a generating set. One sees similarly that $f \circ t$ is an identity map, after checking that elements of the form $(x_1 \otimes \cdots \otimes x_k) \otimes (x_{k+1} \otimes \cdots \otimes x_{k+m})$ generate the R -module $M^{\otimes k} \otimes_R M^{\otimes m}$.

By an entirely analogous proof, one can show that

$$(M_1 \otimes_R \cdots \otimes_R M_k) \otimes_R (M_{k+1} \otimes_R \cdots \otimes_R M_{k+m}) \cong (M_1 \otimes_R \cdots \otimes_R M_{k+m}) \quad (20.12)$$

for any R -modules M_1, \dots, M_{k+m} (Exercise 34). More generally, no matter how we insert parentheses into $M_1 \otimes_R \cdots \otimes_R M_{k+m}$ to indicate nested tensor product constructions, we still obtain a module isomorphic to the original tensor product. Another approach to proving these isomorphisms is to show that both modules solve the same universal mapping problem.

20.10 Tensor Product of Maps

Let M_1, \dots, M_n and P_1, \dots, P_n be R -modules, and let $f_k : M_k \rightarrow P_k$ be an R -linear map for $1 \leq k \leq n$. We now show that there exists a unique R -linear map $g : M_1 \otimes_R \cdots \otimes_R M_n \rightarrow P_1 \otimes_R \cdots \otimes_R P_n$ given on generators by

$$g(x_1 \otimes x_2 \otimes \cdots \otimes x_n) = f_1(x_1) \otimes f_2(x_2) \otimes \cdots \otimes f_n(x_n) \quad (20.13)$$

for all $x_k \in M_k$. The map g is often denoted $f_1 \otimes f_2 \otimes \cdots \otimes f_n$ or $\bigotimes_{k=1}^n f_k : \bigotimes_{k=1}^n M_k \rightarrow \bigotimes_{k=1}^n P_k$ and called the *tensor product of the maps* f_k .

As in previous sections, to obtain the R -linear map g we must invoke the UMP for tensor products. Define a map $h : M_1 \times \cdots \times M_n \rightarrow P_1 \otimes_R \cdots \otimes_R P_n$ by

$$h(x_1, x_2, \dots, x_n) = f_1(x_1) \otimes f_2(x_2) \otimes \cdots \otimes f_n(x_n)$$

for all $x_k \in M_k$. The function h is R -multilinear, since

$$\begin{aligned} h(x_1, \dots, x_k + x'_k, \dots, x_n) &= f_1(x_1) \otimes \cdots \otimes f_k(x_k + x'_k) \otimes \cdots \otimes f_n(x_n) \\ &= f_1(x_1) \otimes \cdots \otimes (f_k(x_k) + f_k(x'_k)) \otimes \cdots \otimes f_n(x_n) \\ &= f_1(x_1) \otimes \cdots \otimes f_k(x_k) \otimes \cdots \otimes f_n(x_n) \\ &\quad + f_1(x_1) \otimes \cdots \otimes f_k(x'_k) \otimes \cdots \otimes f_n(x_n) \\ &= h(x_1, \dots, x_k, \dots, x_n) + h(x_1, \dots, x'_k, \dots, x_n); \end{aligned}$$

and (for all $r \in R$)

$$\begin{aligned} h(x_1, \dots, rx_k, \dots, x_n) &= f_1(x_1) \otimes \cdots \otimes f_k(rx_k) \otimes \cdots \otimes f_n(x_n) \\ &= f_1(x_1) \otimes \cdots \otimes rf_k(x_k) \otimes \cdots \otimes f_n(x_n) \\ &= r(f_1(x_1) \otimes \cdots \otimes f_k(x_k) \otimes \cdots \otimes f_n(x_n)) = rh(x_1, \dots, x_k, \dots, x_n). \end{aligned}$$

Applying the UMP for tensor products to the R -multilinear map h , we obtain a unique R -linear map g satisfying (20.13).

With the same setup as above, suppose further that we have R -modules Q_1, \dots, Q_n and R -linear maps $g_k : P_k \rightarrow Q_k$ for $1 \leq k \leq n$. Then we have R -linear maps $\bigotimes_{k=1}^n f_k : \bigotimes_{k=1}^n M_k \rightarrow \bigotimes_{k=1}^n P_k$, $\bigotimes_{k=1}^n g_k : \bigotimes_{k=1}^n P_k \rightarrow \bigotimes_{k=1}^n Q_k$, and (for each k) $g_k \circ f_k : M_k \rightarrow Q_k$. We claim that

$$\left(\bigotimes_{k=1}^n g_k \right) \circ \left(\bigotimes_{k=1}^n f_k \right) = \bigotimes_{k=1}^n (g_k \circ f_k). \quad (20.14)$$

Both sides are R -linear maps from $\bigotimes_{k=1}^n M_k$ to $\bigotimes_{k=1}^n Q_k$, so it suffices to check that these

functions have the same effect on all generators of $\bigotimes_{k=1}^n M_k$. To check this, note that for all $x_k \in M_k$,

$$\begin{aligned}(g_1 \otimes \cdots \otimes g_k) \circ (f_1 \otimes \cdots \otimes f_k)(x_1 \otimes \cdots \otimes x_k) &= (g_1 \otimes \cdots \otimes g_k)(f_1(x_1) \otimes \cdots \otimes f_k(x_k)) \\ &= g_1(f_1(x_1)) \otimes \cdots \otimes g_k(f_k(x_k)) = ((g_1 \circ f_1) \otimes \cdots \otimes (g_k \circ f_k))(x_1 \otimes \cdots \otimes x_k).\end{aligned}$$

Similarly, we have $\bigotimes_{k=1}^n \text{id}_{M_k} = \text{id}_{\bigotimes_{k=1}^n M_k}$ because both sides are R -linear and

$$\begin{aligned}(\text{id}_{M_1} \otimes \cdots \otimes \text{id}_{M_n})(x_1 \otimes \cdots \otimes x_n) &= \text{id}_{M_1}(x_1) \otimes \cdots \otimes \text{id}_{M_n}(x_n) \\ &= x_1 \otimes \cdots \otimes x_n = \text{id}_{\bigotimes_{k=1}^n M_k}(x_1 \otimes \cdots \otimes x_n).\end{aligned}$$

We can execute a similar construction to induce linear maps on exterior and symmetric powers of R -modules. Suppose $f : M \rightarrow P$ is an R -linear map between R -modules M and P , and $n \in \mathbb{N}$. We claim there exists a unique R -linear map

$$\bigwedge^n f : \bigwedge^n M \rightarrow \bigwedge^n P$$

given on generators by

$$\left(\bigwedge^n f \right) (x_1 \wedge x_2 \wedge \cdots \wedge x_n) = f(x_1) \wedge f(x_2) \wedge \cdots \wedge f(x_n)$$

for all $x_i \in M$. This map is called the n 'th exterior power of f . To obtain this map, define $h : M^n \rightarrow \bigwedge^n P$ by $h(x_1, \dots, x_n) = f(x_1) \wedge \cdots \wedge f(x_n)$. Using R -linearity of f , one checks as above that h is R -multilinear. The map h is alternating as well, since $x_i = x_j$ for $i < j$ gives $f(x_i) = f(x_j)$, hence $f(x_1) \wedge \cdots \wedge f(x_i) \wedge \cdots \wedge f(x_j) \wedge \cdots \wedge f(x_n) = 0$. Applying the UMP for exterior powers to h , we obtain a unique R -linear map $\bigwedge^n f$ satisfying the formula above. If Q is another R -module and $g : P \rightarrow Q$ is another R -linear map, one proves that $(\bigwedge^n g) \circ (\bigwedge^n f) = \bigwedge^n(g \circ f)$ by a calculation on generators analogous to the one used above to prove (20.14). Similarly, one sees that $\bigwedge^n \text{id}_M = \text{id}_{\bigwedge^n M}$.

With the same setup as the previous paragraph, the same method proves the existence of a unique R -linear map

$$\text{Sym}^n f : \text{Sym}^n M \rightarrow \text{Sym}^n P$$

given on generators by

$$(\text{Sym}^n f)(x_1 x_2 \cdots x_n) = f(x_1) f(x_2) \cdots f(x_n)$$

for $x_i \in M$. This map is called the n 'th symmetric power of f . One may check that $(\text{Sym}^n g) \circ (\text{Sym}^n f) = \text{Sym}^n(g \circ f)$ and $\text{Sym}^n \text{id}_M = \text{id}_{\text{Sym}^n M}$.

20.11 Bases and Multilinear Maps

In this section, we study a universal mapping property involving multilinear maps defined on a product of free R -modules. We will see that each such multilinear map is uniquely determined by its effect on the product of bases of the given modules. By comparing this result to the UMP for tensor products, we will obtain an explicit basis for a tensor product of free R -modules.

To begin, assume that M_1, \dots, M_n are free R -modules with respective bases X_1, \dots, X_n . Let $M = M_1 \times \dots \times M_n$, let $X = X_1 \times \dots \times X_n$, and let $i : X \rightarrow M$ be the inclusion map. Then the following universal mapping property holds.

UMP for Multilinear Maps on Free Modules. For any R -module N and any function $f : X \rightarrow N$, there exists a unique R -multilinear map $T : M \rightarrow N$ such that $f = T \circ i$.

$$\begin{array}{ccc} X & \xrightarrow{i} & M \\ & \searrow f & \downarrow T \\ & & N \end{array}$$

Equivalently: for any R -module N , there is a bijection from the set

$$A = \{R\text{-multilinear maps } T : M \rightarrow N\}$$

onto the set

$$B = \{\text{arbitrary maps } f : X \rightarrow N\},$$

which sends each $T \in A$ to $T \circ i \in B$. Informally, we say that each function $f : X \rightarrow N$ “extends by multilinearity” to a unique multilinear map $T : M \rightarrow N$.

To prove the UMP, fix the module N and the map $f : X \rightarrow N$. Suppose $T : M \rightarrow N$ is R -multilinear and extends f . Since each M_k is free with basis X_k , each element $m_k \in M_k$ can be written uniquely as a finite R -linear combination of elements of X_k , say $m_k = \sum_{j_k=1}^{s_k} c_{j_k,k} x_{j_k,k}$ with each $x_{j_k,k} \in X_k$ and each $c_{j_k,k} \in R$. Using (20.4), we see that $T(m_1, m_2, \dots, m_n)$ must be given by the formula

$$\sum_{j_1=1}^{s_1} \sum_{j_2=1}^{s_2} \cdots \sum_{j_n=1}^{s_n} c_{j_1,1} c_{j_2,2} \cdots c_{j_n,n} f(x_{j_1,1}, x_{j_2,2}, \dots, x_{j_n,n}).$$

This shows that T , if it exists at all, is uniquely determined by f . To prove existence, use the previous formula as the definition of $T(m_1, \dots, m_n)$. It is routine to check that T extends f . To see that T is multilinear, fix an index k , a scalar $r \in R$, and $m'_k = \sum_{j_k=1}^{s_k} d_{j_k,k} x_{j_k,k} \in M_k$. On one hand, since $rm_k = \sum_{j_k=1}^{s_k} (rc_{j_k,k}) x_{j_k,k}$, we get

$$\begin{aligned} T(m_1, \dots, rm_k, \dots, r_n) &= \sum_{j_1=1}^{s_1} \cdots \sum_{j_n=1}^{s_n} c_{j_1,1} \cdots (rc_{j_k,k}) \cdots c_{j_n,n} f(x_{j_1,1}, \dots, x_{j_n,n}) \\ &= r \sum_{j_1=1}^{s_1} \cdots \sum_{j_n=1}^{s_n} c_{j_1,1} \cdots c_{j_k,k} \cdots c_{j_n,n} f(x_{j_1,1}, \dots, x_{j_n,n}) \\ &= rT(m_1, \dots, m_k, \dots, r_n), \end{aligned}$$

where the second equality used commutativity of R . On the other hand, since

$m_k + m'_k = \sum_{j_k=1}^{s_k} (c_{j_k,k} + d_{j_k,k})x_{j_k,k}$, we get

$$\begin{aligned}
T(m_1, \dots, m_k + m'_k, \dots, m_n) &= \sum_{j_1=1}^{s_1} \cdots \sum_{j_n=1}^{s_n} c_{j_1,1} \cdots (c_{j_k,k} + d_{j_k,k}) \cdots c_{j_n,n} f(x_{j_1,1}, \dots, x_{j_n,n}) \\
&= \sum_{j_1=1}^{s_1} \cdots \sum_{j_n=1}^{s_n} c_{j_1,1} \cdots c_{j_k,k} \cdots c_{j_n,n} f(x_{j_1,1}, \dots, x_{j_n,n}) \\
&\quad + \sum_{j_1=1}^{s_1} \cdots \sum_{j_n=1}^{s_n} c_{j_1,1} \cdots d_{j_k,k} \cdots c_{j_n,n} f(x_{j_1,1}, \dots, x_{j_n,n}) \\
&= T(m_1, \dots, m_k, \dots, m_n) + T(m_1, \dots, m'_k, \dots, m_n).
\end{aligned}$$

So T is multilinear, completing the proof of the UMP.

20.12 Bases for Tensor Products of Free R -Modules

We continue to assume that M_1, \dots, M_n are free R -modules with respective bases X_1, \dots, X_n . We will now use the UMP proved in §20.11 to construct a new solution to the universal mapping problem for multilinear maps posed in §20.4. Our ultimate goal is to show that $\{x_1 \otimes x_2 \otimes \cdots \otimes x_n : x_k \in X_k\}$ is a basis of $M_1 \otimes_R M_2 \otimes_R \cdots \otimes_R M_n$.

We have already constructed an R -module $N = M_1 \otimes_R \otimes \cdots \otimes_R M_n$ and a multilinear map $j : M_1 \times \cdots \times M_n \rightarrow N$ such that (N, j) solves the UMP posed in §20.4. To build the second solution, let N' be a free R -module with basis $X = X_1 \times \cdots \times X_n$ (see §17.12). Define a multilinear map $j' : M_1 \times \cdots \times M_n \rightarrow N'$ as follows. For each $x = (x_1, \dots, x_n) \in X$, define $j'(x) = x \in N'$. Using the UMP proved in §20.11, j' extends by multilinearity to a unique multilinear map with domain $M_1 \times \cdots \times M_n$. Explicitly, we have

$$j' \left(\sum_{k_1 \geq 1} c_{k_1,1} x_{k_1,1}, \dots, \sum_{k_n \geq 1} c_{k_n,n} x_{k_n,n} \right) = \sum_{k_1 \geq 1} \cdots \sum_{k_n \geq 1} c_{k_1,1} \cdots c_{k_n,n} (x_{k_1,1}, \dots, x_{k_n,n}) \in N'.$$

To prove that (N', j') solves the UMP, suppose P is any R -module and $f : M_1 \times \cdots \times M_n \rightarrow P$ is any multilinear map. We must show there exists a unique R -linear map $g' : N' \rightarrow P$ with $f = g' \circ j'$.

$$\begin{array}{ccc}
M_1 \times \cdots \times M_n & \xrightarrow{j'} & N' \\
f \searrow & & \downarrow g' \\
& & P
\end{array}$$

To see that g' exists, define $g'(x) = f(x)$ for all $x \in X$ and extend g' by linearity to an R -linear map $g' : N' \rightarrow P$ (using the UMP for free R -modules). Now f and $g' \circ j'$ are two multilinear maps from $M_1 \times \cdots \times M_n$ to P that agree on $X = X_1 \times \cdots \times X_n$, since $f(x) = g'(x) = g'(j'(x)) = (g' \circ j')(x)$ for all $x \in X$. By the uniqueness property in the UMP from §20.11, $f = g' \circ j'$ as needed. To see that g' is unique, suppose we also had $f = h' \circ j'$ for some R -linear $h' : N' \rightarrow P$. Then $h'(x) = h'(j'(x)) = f(x) = g'(x)$ for all $x \in X$. Two R -linear maps that agree on a basis of N' must be equal, so $g' = h'$.

Now we know that (N', j') and (N, j) both solve the same UMP. Thus there exists a unique R -module isomorphism $g' : N' \rightarrow N$ such that $j = g' \circ j'$.

$$\begin{array}{ccc} M_1 \times \cdots \times M_n & \xrightarrow{j'} & N' \\ & \searrow j & \downarrow g' \\ & M_1 \otimes_R \cdots \otimes_R M_n & \end{array}$$

The isomorphism g' sends the R -basis $X = X_1 \times \cdots \times X_n$ of N' onto an R -basis of the tensor product $M_1 \otimes_R \cdots \otimes_R M_n$. Explicitly, g' sends $x = (x_1, \dots, x_n) \in X$ to $g'(x) = g'(j'(x)) = j(x) = x_1 \otimes \cdots \otimes x_n$. We have now proved that

$$\{x_1 \otimes x_2 \otimes \cdots \otimes x_n : x_k \in X_k\}$$

is a basis for the R -module $M_1 \otimes_R M_2 \otimes_R \cdots \otimes_R M_n$. In particular, if $\dim(M_k) = d_k < \infty$ for each k , we see that $\dim(\bigotimes_{k=1}^n M_k) = d_1 d_2 \cdots d_n$.

20.13 Bases and Alternating Maps

Let M be a free R -module with basis X . We want to use X to build a basis for the exterior power $\bigwedge^n M$. First, we need to establish a universal mapping property for alternating maps. We fix a total ordering $<$ on X ; for instance, if $X = \{x_1, x_2, \dots, x_m\}$ is finite, we can use the ordering $x_1 < x_2 < \cdots < x_m$. Let $X^n = \{(z_1, \dots, z_n) : z_i \in X\}$, and let

$$X^n_< = \{(z_1, \dots, z_n) \in X^n : z_1 < z_2 < \cdots < z_n\}$$

be the set of strictly increasing sequences of n basis elements. (If $n > |X|$, then $X^n_<$ is empty.) Let $i : X^n_< \rightarrow M^n$ be the inclusion mapping of $X^n_<$ into the product module M^n .

UMP for Alternating Maps on Free Modules. For any R -module N and any function $f : X^n_< \rightarrow N$, there exists a unique alternating map $T : M^n \rightarrow N$ such that $f = T \circ i$.

$$\begin{array}{ccc} X^n_< & \xrightarrow{i} & M^n \\ & \searrow f & \downarrow T \\ & & N \end{array}$$

Equivalently: for any R -module N , there is a bijection from the set

$$A = \{\text{alternating maps } T : M^n \rightarrow N\}$$

onto the set

$$B = \{\text{arbitrary maps } f : X^n_< \rightarrow N\},$$

which sends each $T \in A$ to $T \circ i \in B$. Intuitively, the UMP says that we can build alternating maps by deciding where to send each strictly increasing list of n basis elements, and the alternating map is uniquely determined by these decisions.

To prove the UMP, fix the module N and the map $f : X^n_< \rightarrow N$. We first extend f to a map $g : X^n \rightarrow N$. Given $z = (z_1, z_2, \dots, z_n) \in X^n$, consider two cases. If $z_i = z_j$ for some $i \neq j$, let $g(z) = 0_N$. If all z_i 's are distinct, let $\text{sort}(z) \in X^n_<$ be

the unique sequence obtained by rearranging the entries of z into increasing order. We can write $\text{sort}(z) = (z_{w(1)}, z_{w(2)}, \dots, z_{w(n)})$ for a unique $w \in S_n$. In this case, define $g(z) = \text{sgn}(w)f(\text{sort}(z)) \in N$. Recall from Chapter 2 that $\text{sgn}(w) = (-1)^{\text{inv}(w)}$, where $\text{inv}(w)$ is the number of interchanges of adjacent elements needed to pass from $\text{sort}(z)$ to z or vice versa. It follows from this that if z and z' differ by interchanging two adjacent elements, then $g(z) = -g(z')$. In turn, we deduce that if z and z' differ by interchanging any two elements, then $g(z) = -g(z')$.

From §20.11, we know that $g : X^n \rightarrow N$ extends uniquely by multilinearity to give a multilinear map $T : M^n \rightarrow N$. Since T extends g , T also extends f , so $f = T \circ i$. We must show that T is alternating. Fix $(m_1, \dots, m_n) \in M^n$ with $m_i = m_j$ for some $i \neq j$. Write $m_k = \sum_{z \in X} r(k, z)z$ for some $r(k, z) \in R$. Since $m_i = m_j$, $r(i, z) = r(j, z)$ for all $z \in X$. Now, since T is multilinear,

$$\begin{aligned} T(m_1, \dots, m_n) &= \sum_{z_1 \in X} \cdots \sum_{z_n \in X} r(1, z_1) \cdots r(n, z_n) T(z_1, \dots, z_n) \\ &= \sum_{z=(z_1, \dots, z_n) \in X^n} r(1, z_1) \cdots r(i, z_i) \cdots r(j, z_j) \cdots r(n, z_n) g(z_1, \dots, z_n). \end{aligned}$$

By definition of g , we can drop all terms in this sum in which two entries of z are equal. The other terms can be split into pairs $z = (z_1, \dots, z_i, \dots, z_j, \dots, z_n)$ and $z' = (z_1, \dots, z_j, \dots, z_i, \dots, z_n)$ by switching the basis elements in positions i and j . By the observations above, $g(z') = -g(z)$. Now, the term indexed by z is

$$r(1, z_1) \cdots r(i, z_i) \cdots r(i, z_j) \cdots r(n, z_n) g(z),$$

whereas the term indexed by z' is

$$r(1, z_1) \cdots r(i, z_j) \cdots r(i, z_i) \cdots r(n, z_n) g(z').$$

Since R is commutative, these two terms add together to give zero. Thus, adding up all these pairs, the total effect is that $T(m_1, \dots, m_i, \dots, m_i, \dots, m_n) = 0_N$.

We must still prove uniqueness of T . Suppose $T' : M^n \rightarrow N$ is another alternating map that extends $f : X^n \rightarrow N$. If we can show that T' extends $g : X^n \rightarrow N$, we can conclude that $T = T'$ using the known uniqueness property from the UMP in §20.11. Suppose $z = (z_1, \dots, z_n) \in X^n$ has two repeated entries. Then $T'(z) = 0 = g(z)$ since T' is alternating. Otherwise, when all entries of z are distinct, let $\text{sort}(z) = (z_{w(1)}, \dots, z_{w(n)})$ for some $w \in S_n$. Property (AC3) in §20.2 holds for the alternating map T' , so

$$T'(z) = \text{sgn}(w)T'(\text{sort}(z)) = \text{sgn}(w)f(\text{sort}(z)) = g(z).$$

Thus T and T' both extend g , hence $T = T'$.

20.14 Bases for Exterior Powers of Free Modules

We continue to assume that M is a free R -module with a basis X totally ordered by $<$. Following the pattern of §20.12, we now prove that

$$\{z_1 \wedge z_2 \wedge \cdots \wedge z_n : z_i \in X, z_1 < z_2 < \cdots < z_n\}$$

is a basis of the R -module $\bigwedge^n M$.

We have already constructed an R -module $N = \bigwedge^n M$ and an alternating map $i : M^n \rightarrow N$ such that (N, i) solves the UMP posed in §20.5. We now construct a second solution to this UMP by letting N' be a free R -module with basis $X_{<}^n = \{(z_1, \dots, z_n) \in X^n : z_1 < \dots < z_n\}$. Noting that $X_{<}^n$ is a subset of both M^n and N' , we can define an alternating map $i' : M^n \rightarrow N'$ by letting $i'(z) = z$ for each $z \in X_{<}^n$ and extending i' to M^n by the UMP in §20.13. To prove that (N', i') solves the UMP for exterior powers, suppose P is any R -module and $f : M^n \rightarrow P$ is any alternating map. We must show there exists a unique R -linear map $g' : N' \rightarrow P$ with $f = g' \circ i'$.

$$\begin{array}{ccc} M^n & \xrightarrow{i'} & N' \\ f \searrow & & \downarrow g' \\ & & P \end{array}$$

To see that g' exists, define $g'(x) = f(x)$ for all $x \in X_{<}^n$ and extend g' by linearity to an R -linear map $g' : N' \rightarrow P$ (using the UMP for free R -modules). Now f and $g' \circ i'$ are two alternating maps from M^n to P that agree on $X_{<}^n$, since $f(x) = g'(x) = g'(i'(x)) = (g' \circ i')(x)$ for all $x \in X_{<}^n$. By the uniqueness property in the UMP from §20.13, $f = g' \circ i'$ as needed. To see that g' is unique, suppose we also had $f = h' \circ i'$ for some R -linear $h' : N' \rightarrow P$. Then $h'(x) = h'(i'(x)) = f(x) = g'(x)$ for all $x \in X_{<}^n$. Two R -linear maps that agree on a basis of N' must be equal, so $g' = h'$.

Now we know that (N', i') and (N, i) both solve the same UMP. Thus there exists a unique R -module isomorphism $g' : N' \rightarrow N$ such that $i = g' \circ i'$.

$$\begin{array}{ccc} M^n & \xrightarrow{i'} & N' \\ i \searrow & & \downarrow g' \\ & & \bigwedge^n M \end{array}$$

The isomorphism g' sends the R -basis $X_{<}^n$ of N' onto an R -basis of the exterior power $\bigwedge^n M$. Explicitly, g' sends $z = (z_1, \dots, z_n) \in X_{<}^n$ to $g'(z) = g'(i'(z)) = i(z) = z_1 \wedge \dots \wedge z_n$. We have now proved that

$$i[X_{<}^n] = \{z_1 \wedge z_2 \wedge \dots \wedge z_n : z_i \in X, z_1 < \dots < z_n\}$$

is a basis for the R -module $\bigwedge^n M$. In particular, if $\dim(M) = d < \infty$, we see that $\dim(\bigwedge^n M) = \binom{d}{n} = \frac{d!}{n!(d-n)!}$ for $0 \leq n \leq d$, and $\dim(\bigwedge^n M) = 0$ for $n > d$.

20.15 Bases for Symmetric Powers of Free Modules

We still assume M is a free R -module with totally ordered basis X . Let X_{\leq}^n be the set of all sequences $z = (z_1, \dots, z_n) \in X^n$ with $z_1 \leq z_2 \leq \dots \leq z_n$ relative to the total ordering $<$ on X . Let $i : X_{\leq}^n \rightarrow M^n$ be the inclusion map given by $i(z) = z$ for $z \in X_{\leq}^n$.

UMP for Symmetric Maps on Free Modules. For any R -module N and any function

$f : X_{\leq}^n \rightarrow N$, there exists a unique symmetric map $T : M^n \rightarrow N$ such that $f = T \circ i$.

$$\begin{array}{ccc} X_{\leq}^n & \xrightarrow{i} & M^n \\ f \searrow & & \downarrow T \\ & & N \end{array}$$

Equivalently: for any R -module N , there is a bijection from the set

$$A = \{\text{symmetric maps } T : M^n \rightarrow N\}$$

onto the set

$$B = \{\text{arbitrary maps } f : X_{\leq}^n \rightarrow N\},$$

which sends each $T \in A$ to $T \circ i \in B$. So, we can build symmetric maps uniquely by specifying where to send each weakly increasing list of n basis elements.

The proof is very similar to the one in §20.13, so we leave certain details as exercises. Fix the module N and the map $f : X_{\leq}^n \rightarrow N$. Define $g : X^n \rightarrow N$ by $g(z) = f(\text{sort}(z))$ for all $z \in X^n$, where $\text{sort}(z)$ is the unique weakly increasing sequence of basis elements that can be obtained by sorting the entries of z . From §20.11, we know that $g : X^n \rightarrow N$ extends uniquely by multilinearity to give a multilinear map $T : M^n \rightarrow N$. Since T extends g , T also extends f , so $f = T \circ i$. Using multilinearity and the definition of g , one checks that T is a symmetric map. For uniqueness of T , suppose $T' : M^n \rightarrow N$ is another symmetric map that extends $f : X_{\leq}^n \rightarrow N$. Using symmetry and the definition of g , check that T and T' must both extend $g : X^n \rightarrow N$. Hence, $T = T'$ follows from the known uniqueness property in the UMP from §20.11.

Having proved the UMP, we use it to show that

$$i[X_{\leq}^n] = \{z_1 z_2 \cdots z_n : z_i \in X, z_1 \leq z_2 \leq \cdots \leq z_n\}$$

is a basis of the R -module $\text{Sym}^n M$. We have already constructed an R -module $N = \text{Sym}^n M$ and a symmetric map $i : M^n \rightarrow N$ such that (N, i) solves the UMP posed in §20.6. We now construct a second solution to this UMP by letting N' be a free R -module with basis X_{\leq}^n . Define a symmetric map $i' : M^n \rightarrow N'$ by sending z to z for each $z \in X_{\leq}^n$ and extending to M^n by the UMP just proved. One may now verify, exactly as in the case of exterior powers, that (N', i') solves the same UMP that (N, i) does. So there is a unique isomorphism between N' and N compatible with the universal maps, and this isomorphism maps the R -basis X_{\leq}^n of N' to the claimed basis of $N = \text{Sym}^n M$.

In particular, if \bar{M} is a free module of dimension $d < \infty$, then for all $n \geq 0$, $\dim(\text{Sym}^n M)$ is the number of weakly increasing sequences of length n drawn from a d -letter totally ordered alphabet. A counting argument (Exercise 41) shows that $\dim(\text{Sym}^n M) = \binom{d+n-1}{n}$.

20.16 Tensor Product of Matrices

Suppose M is a free R -module with ordered basis $X = (x_1, \dots, x_m)$ and N is a free R -module with ordered basis $Y = (y_1, \dots, y_n)$. We have seen that the list of mn basic tensors

$$Z = (x_1 \otimes y_1, x_2 \otimes y_1, \dots, x_m \otimes y_1, x_1 \otimes y_2, x_2 \otimes y_2, \dots, x_m \otimes y_2, \dots, x_m \otimes y_n)$$

is an ordered basis of $M \otimes_R N$. Now suppose $f : M \rightarrow M$ and $g : N \rightarrow N$ are R -linear maps. Let $A \in M_m(R)$ be the matrix of f relative to the basis X , and let $B \in M_n(R)$

be the matrix of g relative to the basis Y . What is the matrix C of the R -linear map $f \otimes g : M \otimes_R N \rightarrow M \otimes_R N$ relative to the basis Z ?

To answer this question, first recall that $f(x_j) = \sum_{i=1}^m A(i, j)x_i$ for all $j \in [m] = \{1, 2, \dots, m\}$, and $g(y_j) = \sum_{i=1}^n B(i, j)y_i$ for all $j \in [n]$. Let us label the rows and columns of the matrix C with the elements of Z in the order they appear above. To compute the entries in the column of C labeled by $x_i \otimes y_j$, we apply $f \otimes g$ to this element, obtaining

$$\begin{aligned} (f \otimes g)(x_i \otimes y_j) &= f(x_i) \otimes g(y_j) = \left(\sum_{k=1}^m A(k, i)x_k \right) \otimes \left(\sum_{\ell=1}^n B(\ell, j)y_\ell \right) \\ &= \sum_{k=1}^m \sum_{\ell=1}^n A(k, i)B(\ell, j)(x_k \otimes y_\ell). \end{aligned}$$

So, the entry of C in the row labeled $(x_k \otimes y_\ell)$ and the column labeled $(x_i \otimes y_j)$ is $A(k, i)B(\ell, j)$. We write $C = A \otimes B$ and call C the tensor product of the matrices A and B .

Note that C is an $mn \times mn$ matrix, where each entry is a scalar in R . We can also think of C as an $n \times n$ “block matrix” where each block is itself an $m \times m$ matrix. For $1 \leq i, j \leq n$, the i, j -block of C has rows labeled $x_1 \otimes y_i, x_2 \otimes y_i, \dots, x_m \otimes y_i$ and columns labeled $x_1 \otimes y_j, x_2 \otimes y_j, \dots, x_m \otimes y_j$. The above calculation shows that for $1 \leq r, s \leq m$, the r, s -entry of the i, j -block of C is $A(r, s)B(i, j)$. Thus, the entire i, j -block of C is found by multiplying the entire matrix A by the scalar $B(i, j)$. Pictorially, the block matrix C is

$$C = A \otimes B = \begin{bmatrix} B(1, 1)A & B(1, 2)A & \cdots & B(1, n)A \\ B(2, 1)A & B(2, 2)A & \cdots & B(2, n)A \\ \vdots & \vdots & \vdots & \vdots \\ B(n, 1)A & B(n, 2)A & \cdots & B(n, n)A \end{bmatrix}.$$

Suppose now that $f_1 : M \rightarrow M$ and $g_1 : N \rightarrow N$ are also R -linear maps, represented (relative to the bases X and Y) by matrices A_1 and B_1 . We have seen that $(f_1 \circ f) \otimes (g_1 \circ g) = (f_1 \otimes g_1) \circ (f \otimes g)$. Passing to matrices, this yields the matrix identity

$$(A_1 A) \otimes (B_1 B) = (A_1 \otimes B_1)(A \otimes B).$$

Similarly, one readily checks that $(f_1 + f) \otimes g = (f_1 \otimes g) + (f \otimes g)$ and $(rf) \otimes g = r(f \otimes g)$ for all $r \in R$. We therefore see that

$$(A_1 + A) \otimes B = (A_1 \otimes B) + (A \otimes B), \quad (rA) \otimes B = r(A \otimes B),$$

and similar identities hold in the second argument.

20.17 Determinants and Exterior Powers

Let M be a free R -module with ordered basis $X = (x_1, \dots, x_n)$. Given an R -linear map $f : M \rightarrow M$, let $A \in M_n(R)$ be the matrix of f relative to X . For each k with $0 \leq k \leq n$, there is an induced R -linear map $\bigwedge^k f : \bigwedge^k M \rightarrow \bigwedge^k M$. What is the matrix of this map relative to the basis $i[X^n_<]$ of $\bigwedge^k M$?

Before answering this question in general, consider the special case $k = n$. Here, $\bigwedge^n M$ is a one-dimensional free R -module, since its basis $i[X^n_<]$ consists of the single wedge product $x^* = x_1 \wedge x_2 \wedge \dots \wedge x_n$ in which all elements of X appear in order. So the matrix of $\bigwedge^n f$ has

just one entry; we find it by applying $\bigwedge^n f$ to x^* and seeing which multiple of x^* results. Recall $f(x_i) = \sum_{k=1}^n A(k, i)x_k$ for $1 \leq i \leq n$. Using this, we compute:

$$\begin{aligned} \left(\bigwedge^n f\right)(x^*) &= f(x_1) \wedge f(x_2) \wedge \cdots \wedge f(x_n) \\ &= \left(\sum_{k_1=1}^n A(k_1, 1)x_{k_1}\right) \wedge \left(\sum_{k_2=1}^n A(k_2, 2)x_{k_2}\right) \wedge \cdots \wedge \left(\sum_{k_n=1}^n A(k_n, n)x_{k_n}\right) \\ &= \sum_{k_1=1}^n \sum_{k_2=1}^n \cdots \sum_{k_n=1}^n A(k_1, 1)A(k_2, 2) \cdots A(k_n, n)(x_{k_1} \wedge x_{k_2} \wedge \cdots \wedge x_{k_n}). \end{aligned}$$

The last step used the multilinearity of wedge products. To continue simplifying, recall that any wedge product with a repeated term is zero. So instead of summing over all sequences $(k_1, \dots, k_n) \in [n]^n$, it suffices to sum just over the permutations $k = (k_1, \dots, k_n) \in S_n$. Then, using the last identity in §20.5, we obtain

$$\left(\bigwedge^n f\right)(x^*) = \left(\sum_{k \in S_n} \text{sgn}(k) A(k_1, 1)A(k_2, 2) \cdots A(k_n, n)\right) (x_1 \wedge x_2 \wedge \cdots \wedge x_n) = \det(A)x^*.$$

We see that $\bigwedge^n f$ has matrix $[\det(A)]$ relative to the basis (x^*) of $\bigwedge^n M$. This computation shows how the mysterious definition of $\det(A)$ (first given in Chapter 5, equation (5.1)) arises naturally in the theory of exterior powers.

We can now give a one-sentence proof of the product formula for determinants (see §5.13). Take maps $f, g : M \rightarrow M$ whose matrices relative to X are A and B , respectively; since $\bigwedge^n(f \circ g) = (\bigwedge^n f) \circ (\bigwedge^n g)$, applying the preceding result to the maps $f \circ g$, f , and g gives $\det(AB) = \det(A)\det(B)$.

Now we return to the question of computing the matrix C of $\bigwedge^k f$ relative to the basis $i[X^k_<]$, for any k between 0 and n . Each element of $i[X^k_<]$ has the form $x_{i_1} \wedge x_{i_2} \wedge \cdots \wedge x_{i_k}$ for a unique k -element subset $I = \{i_1 < i_2 < \cdots < i_k\}$ of $[n] = \{1, 2, \dots, n\}$. Let us label the rows and columns of C with these subsets. To find the entries of C in the column labeled $J = \{j_1 < j_2 < \cdots < j_k\}$, apply $\bigwedge^k f$ to $x_{j_1} \wedge x_{j_2} \wedge \cdots \wedge x_{j_k}$. We obtain

$$\begin{aligned} f(x_{j_1}) \wedge f(x_{j_2}) \wedge \cdots \wedge f(x_{j_k}) &= \left(\sum_{w_1=1}^n A(w_1, j_1)x_{w_1}\right) \wedge \left(\sum_{w_2=1}^n A(w_2, j_2)x_{w_2}\right) \wedge \cdots \wedge \left(\sum_{w_k=1}^n A(w_k, j_k)x_{w_k}\right) \\ &= \sum_{w_1=1}^n \cdots \sum_{w_k=1}^n A(w_1, j_1) \cdots A(w_k, j_k)x_{w_1} \wedge \cdots \wedge x_{w_k}. \end{aligned}$$

As before, we can discard zero terms to reduce to a sum over $w = (w_1, \dots, w_k)$ with all entries distinct. For each k -element subset $I = \{i_1 < i_2 < \cdots < i_k\}$, there will be $k!$ terms in the sum arising from words w that are rearrangements of the entries of I . If a permutation $f \in S_k$ rearranges I into w , then

$$A(w_1, j_1) \cdots A(w_k, j_k)x_{w_1} \wedge \cdots \wedge x_{w_k} = \text{sgn}(f)A(i_{f(1)}, j_1) \cdots A(i_{f(k)}, j_k)x_{i_1} \wedge \cdots \wedge x_{i_k}.$$

These $k!$ terms in the sum, and no others, will contribute to the coefficient of $x_{i_1} \wedge \cdots \wedge x_{i_k}$. It follows that *the entry of C in the row labeled I and the column labeled J is*

$$\sum_{f \in S_k} \text{sgn}(f)A(i_{f(1)}, j_1) \cdots A(i_{f(k)}, j_k) = \det(A_{I,J}),$$

where $\det(A_{I,J})$ denotes the determinant of the $k \times k$ submatrix of A obtained by keeping only the k rows in I and the k columns in J . Thus, the entries in the matrix for $\bigwedge^k f$ are precisely the k 'th order minors of A , which we studied in §18.11.

20.18 From Modules to Algebras

Recall that an *R-algebra* is a ring $(A, +, \star)$ that is also an *R-module*, such that $c(x \star y) = (cx) \star y = x \star (cy)$ for all $x, y \in A$ and all $c \in R$. An *R-algebra homomorphism* is a map between *R-algebras* that is both a ring homomorphism and an *R-linear map*. Given any *R-module* M , our goal is to build an *R-algebra* $T(M)$, called the *tensor algebra of M*, solving the following universal mapping problem.

UMP for Tensor Algebras. Given an *R-module* M , construct an *R-algebra* $T(M)$ and an *R-linear map* $i : M \rightarrow T(M)$ such that, for any *R-algebra* A and any *R-linear map* $f : M \rightarrow A$, there exists a unique *R-algebra homomorphism* $g : T(M) \rightarrow A$ with $f = g \circ i$.

$$\begin{array}{ccc} M & \xrightarrow{i} & T(M) \\ & \searrow f & \downarrow g \\ & & A \end{array}$$

Construction of $T(M)$. Let $T(M)$ be the (external) direct sum $\bigoplus_{k=0}^{\infty} M^{\otimes k}$, where $M^{\otimes 0} = R$. By definition, an element of $T(M)$ is an infinite sequence $z = (z_0, z_1, \dots, z_k, \dots)$, where $z_k \in M^{\otimes k}$ and all but finitely many z_k 's are zero. We already know $T(M)$ is an *R-module*, and (using Exercise 12) there is an injective *R-linear map* $i : M \rightarrow T(M)$ given by $i(x) = (0, x, 0, 0, \dots)$ for $x \in M$. We must define an algebra structure on $T(M)$ by specifying the ring multiplication $\star : T(M) \times T(M) \rightarrow T(M)$.

In §20.9, we constructed (for each $m, k \in \mathbb{N}^+$) an *R-linear map* $\nu_{k,m} : M^{\otimes k} \otimes_R M^{\otimes m} \rightarrow M^{\otimes(k+m)}$ given on generators by

$$\nu_{k,m}((x_1 \otimes \cdots \otimes x_k) \otimes (x_{k+1} \otimes \cdots \otimes x_{k+m})) = x_1 \otimes \cdots \otimes x_k \otimes x_{k+1} \otimes \cdots \otimes x_{k+m}.$$

Recall that $\nu_{k,m}$ arises (via the UMP for tensor products) from an *R-bilinear map* $\mu_{k,m} : M^{\otimes k} \times M^{\otimes m} \rightarrow M^{\otimes(k+m)}$, which acts on generators by

$$\mu_{k,m}(x_1 \otimes \cdots \otimes x_k, x_{k+1} \otimes \cdots \otimes x_{k+m}) = x_1 \otimes \cdots \otimes x_k \otimes x_{k+1} \otimes \cdots \otimes x_{k+m}.$$

Similarly, we have *R-bilinear maps* $\mu_{0,m}$ and $\mu_{k,0}$ such that (for $r \in R$ and $x_i \in M$)

$$\mu_{0,m}(r, x_1 \otimes \cdots \otimes x_m) = (rx_1) \otimes \cdots \otimes x_m; \quad \mu_{k,0}(x_1 \otimes \cdots \otimes x_k, r) = (rx_1) \otimes \cdots \otimes x_k.$$

We can assemble all the maps $\mu_{k,m}$ to obtain a map $\star : T(M) \times T(M) \rightarrow T(M)$, as follows. Given $y = (y_k : k \geq 0)$ and $z = (z_k : k \geq 0)$ in $T(M)$, define

$$y \star z = \left(\sum_{k+m=n} \mu_{k,m}(y_k, z_m) : n \geq 0 \right).$$

It is tedious but routine to confirm that $(T(M), +, \star)$ satisfies the axioms for a ring and

R -algebra. In particular, the left and right distributive laws for \star follow from the preceding definition and bilinearity of each $\mu_{k,m}$. Associativity of \star follows from the fact that

$$\begin{aligned}\mu_{k+m,p}(\mu_{k,m}(x_1 \otimes \cdots \otimes x_k, y_1 \otimes \cdots \otimes y_m), z_1 \otimes \cdots \otimes z_p) \\ = \mu_{k,m+p}(x_1 \otimes \cdots \otimes x_k, \mu_{m,p}(y_1 \otimes \cdots \otimes y_m, z_1 \otimes \cdots \otimes z_p)),\end{aligned}$$

which holds since both sides equal

$$x_1 \otimes \cdots \otimes x_k \otimes y_1 \otimes \cdots \otimes y_m \otimes z_1 \otimes \cdots \otimes z_p.$$

The multiplicative identity of $T(M)$ is $(1_R, 0, 0, \dots)$.

Now let us verify that $(T(M), i)$ solves the UMP. Let $(A, +, *)$ be any R -algebra, and let $f : M \rightarrow A$ be an R -linear map. We need to build an R -algebra map $g : T(M) \rightarrow A$ with $f = g \circ i$. Define $g_0 : R \rightarrow A$ by $g_0(r) = r \cdot 1_A$ for all $r \in R$. For each $k > 0$, define $f_k : M^k \rightarrow A$ by $f_k(x_1, x_2, \dots, x_k) = f(x_1) * f(x_2) * \cdots * f(x_k)$ for all $x_j \in M$. Since A is an R -algebra and f is R -linear, f_k is k -linear. So we get an induced R -linear map $g_k : M^{\otimes k} \rightarrow A$ given on generators by

$$g_k(x_1 \otimes x_2 \otimes \cdots \otimes x_k) = f(x_1) * f(x_2) * \cdots * f(x_k).$$

Finally, the UMP for direct sums (§19.6) combines all the maps g_k to give an R -linear map $g : T(M) \rightarrow A$ such that $g(z_0, z_1, \dots, z_k, \dots) = \sum_{k \geq 0} g_k(z_k)$. Note that $g \circ i = f$, since for all $z_1 \in M$, $g(i(z_1)) = g(0, z_1, 0, \dots) = g_1(z_1) = f(z_1)$. Since g is already known to be R -linear, we need only check that $g(y \star z) = g(y) * g(z)$ for all $y, z \in T(M)$. Using the definition of \star and the distributive laws, this verification reduces to the fact that $g_{k+m}(\mu_{k,m}(y, z)) = g_k(y) * g_m(z)$ for all $y \in M^{\otimes k}$ and all $z \in M^{\otimes m}$. In turn, this fact holds because

$$\begin{aligned}g_{k+m}(\mu_{k,m}(y_1 \otimes \cdots \otimes y_k, z_1 \otimes \cdots \otimes z_m)) &= f(y_1) * f(y_2) * \cdots * f(y_k) * f(z_1) * \cdots * f(z_m) \\ &= g_k(y_1 \otimes \cdots \otimes y_k) * g_m(z_1 \otimes \cdots \otimes z_m),\end{aligned}$$

$\mu_{k,m}$ is R -bilinear, and g_k, g_m, g_{k+m} are R -linear.

Finally, we must show that g is unique. Suppose $h : T(M) \rightarrow A$ is also an R -algebra homomorphism with $h \circ i = f$. By the uniqueness property in the UMP for direct sums, it suffices to check that $h \circ j_k = g \circ j_k$ for all $k \geq 0$, where j_k is the injection of $M^{\otimes k}$ into $T(M) = \bigoplus_{s \geq 0} M^{\otimes s}$. When $k = 0$, we have

$$h \circ j_0(r) = rh(1_{T(M)}) = r1_A = rg_0(1_{T(M)}) = g \circ j_0(r)$$

for all $r \in R$. When $k = 1$, we have

$$h \circ j_1(x) = h \circ i(x) = f(x) = g \circ i(x) = g \circ j_1(x)$$

for all $x \in M$. When $k \geq 2$, note that $x_1 \otimes x_2 \otimes \cdots \otimes x_k = x_1 \star x_2 \star \cdots \star x_k$ for all $x_j \in M$. Since g and h preserve multiplication and coincide with f on M , we get

$$\begin{aligned}h \circ j_k(x_1 \otimes \cdots \otimes x_k) &= h(x_1 \star \cdots \star x_k) = f(x_1) * f(x_2) * \cdots * f(x_k) \\ &= g(x_1 \star \cdots \star x_k) = g \circ j_k(x_1 \otimes \cdots \otimes x_k).\end{aligned}$$

Thus $g = h$, completing the proof of the uniqueness assertion in the UMP.

In general, the algebra $T(M)$ is not commutative. By replacing tensor powers of M by symmetric powers of M throughout the preceding construction, one can build a commutative algebra $\text{Sym}(M)$ such that any R -linear map from M into a commutative algebra A uniquely extends to an R -algebra map from $\text{Sym}(M)$ to A (Exercise 51). $\text{Sym}(M)$ is called the *symmetric algebra of M* . A similar construction using exterior powers produces an algebra $\Lambda(M)$ solving an appropriate UMP (Exercise 52). $\Lambda(M)$ is called the *exterior algebra of M* .

20.19 Summary

In this summary, assume R is a commutative ring and all other capital letters are R -modules unless otherwise stated.

1. *Special Maps.* A map $f : M \rightarrow N$ is R -linear iff $f(m + m') = f(m) + f(m')$ and $f(rm) = rf(m)$ for all $m, m' \in M$ and $r \in R$. A map $f : M_1 \times \cdots \times M_n \rightarrow N$ is R -multilinear iff for $1 \leq i \leq n$, f is R -linear in position i when all other arguments are held fixed. An R -multilinear map $f : M^n \rightarrow N$ is *alternating* iff f has value zero when any two arguments are equal. An R -multilinear map $f : M^n \rightarrow N$ is *anti-commutative* iff $f(m_{w(1)}, \dots, m_{w(n)}) = \text{sgn}(w)f(m_1, \dots, m_n)$ for all $m_k \in M$ and all $w \in S_n$. To prove anti-commutativity, it suffices to check that f changes sign when any two adjacent arguments are switched. Alternating maps are anti-commutative, but the converse only holds when $1_R + 1_R$ is not a zero divisor. An R -multilinear map $f : M^n \rightarrow N$ is *symmetric* iff $f(m_{w(1)}, \dots, m_{w(n)}) = f(m_1, \dots, m_n)$ for all $m_k \in M$ and all $w \in S_n$. To prove symmetry, it suffices to check that f is unchanged when any two adjacent arguments are switched.
2. *Generators for $\bigotimes_{k=1}^n M_k$, $\bigwedge^n M$, and $\text{Sym}^n M$.* Every element of $M_1 \otimes_R \cdots \otimes_R M_n$ is a finite sum of basic tensors $x_1 \otimes \cdots \otimes x_n$ with $x_k \in M_k$. Every element of $\bigwedge^n M$ is a finite sum of elements $x_1 \wedge \cdots \wedge x_n$ with $x_k \in M$. Every element of $\text{Sym}^n M$ is a finite sum of elements $x_1 \cdots x_n$ with $x_k \in M$. If X_k generates M_k , then $\{x_1 \otimes \cdots \otimes x_n : x_k \in X_k\}$ generates $\bigotimes_{k=1}^n M_k$; similarly for $\bigwedge^n M$ and $\text{Sym}^n M$.
3. *Bases for $\bigotimes_{k=1}^n M_k$, $\bigwedge^n M$, and $\text{Sym}^n M$.* If each M_k is free with basis X_k , then $\bigotimes_{k=1}^n M_k$ is free with basis $\{x_1 \otimes \cdots \otimes x_n : x_k \in X_k\}$. If M is free with basis X totally ordered by $<$, then $\bigwedge^n M$ is free with basis $\{x_1 \wedge \cdots \wedge x_n : x_k \in X, x_1 < \cdots < x_n\}$, and $\text{Sym}^n M$ is free with basis $\{x_1 \cdots x_n : x_k \in X, x_1 \leq \cdots \leq x_n\}$. In the case where $\dim(M_k) = d_k < \infty$ and $\dim(M) = d < \infty$, we have $\dim(M_1 \otimes_R \cdots \otimes_R M_n) = d_1 d_2 \cdots d_n$, $\dim(\bigwedge^n M) = \binom{d}{n}$ for $0 \leq n \leq d$ (and is zero otherwise), and $\dim(\text{Sym}^n M) = \binom{d+n-1}{n}$.
4. *Myths about Tensor Products.* It is false that every element of $M_1 \otimes_R \cdots \otimes_R M_n$ must have the form $x_1 \otimes \cdots \otimes x_n$ for some $x_k \in M_k$. It is false that $\{x_1 \otimes \cdots \otimes x_n : x_k \in M_k\}$ must be a basis for $M_1 \otimes_R \cdots \otimes_R M_n$. It is false that we can define a function on $M_1 \otimes_R \cdots \otimes_R M_n$ (with no further work) by specifying where the function sends basic tensors. It is false that a tensor product of nonzero modules must be nonzero.
5. *Tensor Product Isomorphisms.* Five R -module isomorphisms (defined on generators) are:
 - (a) *Commutativity:* $M \otimes_R N \cong N \otimes_R M$ via the map $m \otimes n \mapsto n \otimes m$.
 - (b) *Associativity:* $(M \otimes_R N) \otimes_R P \cong M \otimes_R (N \otimes_R P)$ via the map $(m \otimes n) \otimes p \mapsto m \otimes (n \otimes p)$.
 - (c) *Left Distributivity:* $(M \oplus N) \otimes_R P \cong (M \otimes_R P) \oplus (N \otimes_R P)$ via the map $(m, n) \otimes p \mapsto (m \otimes p, n \otimes p)$.
 - (d) *Right Distributivity:* $P \otimes_R (M \oplus N) \cong (P \otimes_R M) \oplus (P \otimes_R N)$ via the map $p \otimes (m, n) \mapsto (p \otimes m, p \otimes n)$.

- (e) *Identity for \otimes_R :* $R \otimes_R M \cong M$ via the map $r \otimes m \mapsto rm$; similarly, $M \otimes_R R \cong M$.

More generally:

$$M_{w(1)} \otimes_R \cdots \otimes_R M_{w(n)} \cong M_1 \otimes_R \cdots \otimes_R M_n \text{ for all } w \in S_n;$$

$$(M_1 \otimes_R \cdots \otimes_R M_k) \otimes_R (M_{k+1} \otimes_R \cdots \otimes_R M_n) \cong M_1 \otimes_R \cdots \otimes_R M_n;$$

$$\left(\bigoplus_{i_1 \in I_1} M_{i_1} \right) \otimes_R \cdots \otimes_R \left(\bigoplus_{i_n \in I_n} M_{i_n} \right) \cong \bigoplus_{(i_1, \dots, i_n) \in I_1 \times \cdots \times I_n} M_{i_1} \otimes_R \cdots \otimes_R M_{i_n};$$

and all factors of R can be deleted in a tensor product to give an isomorphic tensor product.

6. *Tensor Product of Linear Maps.* Given R -linear maps $f_k : M_k \rightarrow N_k$, there is an induced R -linear map $\bigotimes_{k=1}^n f_k : \bigotimes_{k=1}^n M_k \rightarrow \bigotimes_{k=1}^n N_k$ given on generators by $(f_1 \otimes \cdots \otimes f_n)(x_1 \otimes \cdots \otimes x_n) = f_1(x_1) \otimes \cdots \otimes f_n(x_n)$. Given R -linear maps $g_k : N_k \rightarrow P_k$, we have $\bigotimes_{k=1}^n (g_k \circ f_k) = (\bigotimes_{k=1}^n g_k) \circ (\bigotimes_{k=1}^n f_k)$.
7. *Exterior Powers and Symmetric Powers of Linear Maps.* Given an R -linear map $f : M \rightarrow N$, there are induced R -linear maps $\Lambda^n f : \Lambda^n M \rightarrow \Lambda^n N$ and $\text{Sym}^n f : \text{Sym}^n M \rightarrow \text{Sym}^n N$ that act on generators by

$$(\Lambda^n f)(x_1 \wedge \cdots \wedge x_n) = f(x_1) \wedge \cdots \wedge f(x_n),$$

$$(\text{Sym}^n f)(x_1 \cdots x_n) = f(x_1) \cdots f(x_n).$$

For an R -linear map $g : N \rightarrow P$, we have $\Lambda^n(g \circ f) = (\Lambda^n g) \circ (\Lambda^n f)$ and $\text{Sym}^n(g \circ f) = (\text{Sym}^n g) \circ (\text{Sym}^n f)$.

8. *Tensor Product of Matrices.* Given $A \in M_m(R)$ and $B \in M_n(R)$, $A \otimes B \in M_{mn}(R)$ is the $n \times n$ block matrix whose i, j -block is $B(i, j)A$. If A is the matrix of f and B is the matrix of g relative to certain bases, $A \otimes B$ is the matrix of $f \otimes g$ relative to the tensor product of these bases. For $C \in M_m(R)$ and $D \in M_n(R)$ and $r \in R$, $(CA) \otimes (DB) = (C \otimes D)(A \otimes B)$, $(C + A) \otimes B = (C \otimes B) + (A \otimes B)$, $A \otimes (B + D) = (A \otimes B) + (A \otimes D)$, and $(rA) \otimes B = r(A \otimes B) = A \otimes (rB)$.
9. *Exterior Powers and Determinants.* Suppose X is an n -element basis for M and $f : M \rightarrow M$ is R -linear. If A is the matrix of f relative to X , then $[\det(A)]$ is the matrix of $\Lambda^n f$ relative to the basis $i[X^n]$. Similarly, the matrix of $\Lambda^k f$ relative to the basis $i[X^k]$ has entries $\det(A_{I,J})$, where $I, J \subseteq [n]$ have size k .

Summary of Universal Mapping Properties

1. *UMP for Tensor Products.* Let $j : M_1 \times \cdots \times M_n \rightarrow M_1 \otimes_R \cdots \otimes_R M_n$ send (m_1, \dots, m_n) to $m_1 \otimes \cdots \otimes m_n$. For every R -module P and every R -multilinear map $f : M_1 \times \cdots \times M_n \rightarrow P$, there exists a unique R -linear map $g : M_1 \otimes_R \cdots \otimes_R M_n \rightarrow P$ with $f = g \circ j$.

$$\begin{array}{ccc} M_1 \times \cdots \times M_n & \xrightarrow{j} & M_1 \otimes_R \cdots \otimes_R M_n \\ & \searrow f & \downarrow g \\ & & P \end{array}$$

So *tensor products convert R -multilinear maps to R -linear maps*.

2. *UMP for Exterior Powers.* Let $i : M^n \rightarrow \bigwedge^n M$ send (m_1, \dots, m_n) to $m_1 \wedge \dots \wedge m_n$. For every R -module P and every alternating map $f : M^n \rightarrow P$, there exists a unique R -linear map $g : \bigwedge^n M \rightarrow P$ with $f = g \circ i$.

$$\begin{array}{ccc} M^n & \xrightarrow{i} & \bigwedge^n M \\ & \searrow f & \downarrow g \\ & & P \end{array}$$

So exterior powers convert alternating maps to R -linear maps.

3. *UMP for Symmetric Powers.* Let $i : M^n \rightarrow \text{Sym}^n M$ send (m_1, \dots, m_n) to $m_1 \cdots m_n$. For every R -module P and every symmetric map $f : M^n \rightarrow P$, there exists a unique R -linear map $g : \text{Sym}^n M \rightarrow P$ with $f = g \circ i$.

$$\begin{array}{ccc} M^n & \xrightarrow{i} & \text{Sym}^n M \\ & \searrow f & \downarrow g \\ & & P \end{array}$$

So symmetric powers convert symmetric maps to R -linear maps.

4. *UMP for Multilinear Maps on Free Modules.* Let M_1, \dots, M_n be free R -modules with respective bases X_1, \dots, X_n . Let $i : X_1 \times \dots \times X_n \rightarrow M_1 \times \dots \times M_n$ be the inclusion map. For any R -module N and any function $f : X_1 \times \dots \times X_n \rightarrow N$, there exists a unique R -multilinear map $T : M_1 \times \dots \times M_n \rightarrow N$ such that $f = T \circ i$.

$$\begin{array}{ccc} X_1 \times \dots \times X_n & \xrightarrow{i} & M_1 \times \dots \times M_n \\ & \searrow f & \downarrow T \\ & & N \end{array}$$

So arbitrary functions on the product of bases extend uniquely to multilinear maps.

5. *UMP for Alternating Maps on Free Modules.* Let M be a free R -module with a basis X totally ordered by $<$. Let $X_{<}^n$ be the set of strictly increasing sequences of n elements of X , and let $i : X_{<}^n \rightarrow M^n$ be the inclusion map. For any R -module N and any function $f : X_{<}^n \rightarrow N$, there exists a unique alternating map $T : M^n \rightarrow N$ such that $f = T \circ i$.

$$\begin{array}{ccc} X_{<}^n & \xrightarrow{i} & M^n \\ & \searrow f & \downarrow T \\ & & N \end{array}$$

So arbitrary functions on $X_{<}^n$ extend uniquely to alternating maps.

6. *UMP for Symmetric Maps on Free Modules.* Let M be a free R -module with a basis X totally ordered by $<$. Let X_{\leq}^n be the set of weakly increasing sequences of n elements of X , and let $i : X_{\leq}^n \rightarrow M^n$ be the inclusion map. For any R -module N and any function $f : X_{\leq}^n \rightarrow N$, there exists a unique symmetric map

$T : M^n \rightarrow N$ such that $f = T \circ i$.

$$\begin{array}{ccc} X_{\leq}^n & \xrightarrow{i} & M^n \\ & \searrow f & \downarrow T \\ & & N \end{array}$$

So arbitrary functions on X_{\leq}^n extend uniquely to symmetric maps.

7. *UMP for Tensor Algebras.* For every R -module M , there is an R -algebra $T(M)$ and an R -linear injection $i : M \rightarrow T(M)$ such that for every R -algebra A and every R -linear map $f : M \rightarrow A$, there exists a unique R -algebra homomorphism $g : T(M) \rightarrow A$ with $f = g \circ i$.

$$\begin{array}{ccc} M & \xrightarrow{i} & T(M) \\ & \searrow f & \downarrow g \\ & & A \end{array}$$

So tensor algebras convert R -linear maps (with codomain an R -algebra) into R -algebra maps.

8. *UMP for Symmetric Algebras.* For every R -module M , there is a commutative R -algebra $\text{Sym}(M)$ and an R -linear injection $i : M \rightarrow \text{Sym}(M)$ such that for every commutative R -algebra A and every R -linear map $f : M \rightarrow A$, there exists a unique R -algebra homomorphism $g : \text{Sym}(M) \rightarrow A$ with $f = g \circ i$.

$$\begin{array}{ccc} M & \xrightarrow{i} & \text{Sym}(M) \\ & \searrow f & \downarrow g \\ & & A \end{array}$$

So symmetric algebras convert R -linear maps (with codomain a commutative R -algebra) into R -algebra maps.

9. *UMP for Exterior Algebras.* For every R -module M , there is an R -algebra $\bigwedge(M)$ with $z \star z = 0$ for all $z \in \bigwedge(M)$ and an R -linear injection $i : M \rightarrow \bigwedge(M)$ such that for every R -algebra A and every R -linear map $f : M \rightarrow A$ such that $f(x) * f(x) = 0$ for all $x \in M$, there exists a unique R -algebra homomorphism $g : \bigwedge(M) \rightarrow A$ with $f = g \circ i$.

$$\begin{array}{ccc} M & \xrightarrow{i} & \bigwedge(M) \\ & \searrow f & \downarrow g \\ & & A \end{array}$$

So exterior algebras convert R -linear maps (where all images square to zero) into R -algebra maps.

20.20 Exercises

Unless otherwise stated, assume R is a commutative ring and M, N, P, M_i are R -modules in these exercises.

1. Define $f : R^n \rightarrow R$ by $f(r_1, r_2, \dots, r_n) = r_1 r_2 \cdots r_n$. Prove that f is a symmetric map. Justify each step using the ring axioms.
2. Decide (with proof) whether each map below is \mathbb{R} -linear, \mathbb{R} -bilinear, or neither. For the bilinear maps, say whether the map is symmetric or alternating.
 - (a) $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ given by $f(x, y) = x + y$.
 - (b) $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ given by $f(x, y) = 5xy$.
 - (c) $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ given by $f(x, y) = x^2 - y^2$.
 - (d) $f : M_3(\mathbb{R}) \times M_3(\mathbb{R}) \rightarrow M_3(\mathbb{R})$ given by $f(A, B) = AB^T$.
 - (e) $f : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ given by $f((a, b), (c, d)) = ad - bc$ for $a, b, c, d \in \mathbb{R}$.
 - (f) $I(g, h) = \int_0^1 tg(t)h(t) dt$ where $g, h : [0, 1] \rightarrow \mathbb{R}$ are continuous functions.
 - (g) $C(g, h) = g \circ h - h \circ g$, where $g, h : \mathbb{R}^n \rightarrow \mathbb{R}^n$ are linear maps.
3. Define $f : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ by $f(v, w) = v^T Aw$, where $A \in M_n(\mathbb{R})$ is a fixed matrix and v, w are column vectors.
 - (a) Prove that f is \mathbb{R} -bilinear.
 - (b) Under what conditions on A is f a linear map?
 - (c) Under what conditions on A is f an alternating map?
 - (d) Under what conditions on A is f a symmetric map?
4. Prove that every \mathbb{R} -bilinear map $f : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ has the form $f(v, w) = v^T Aw$ for some $A \in M_n(\mathbb{R})$.
5. (a) Prove (20.3) by induction on s . (b) Deduce (20.4) from (20.3).
6. Give an example of an anti-commutative bilinear map that is not alternating.
7. (a) Give an example of a nonzero bilinear map that is symmetric and alternating.
 - (b) Find conditions on n and R under which every symmetric alternating map $f : M^n \rightarrow P$ must be zero.
8. (a) Prove: for all multilinear $f : M_1 \times \cdots \times M_n \rightarrow P$, if $m_i = 0$ for some i then $f(m_1, \dots, m_n) = 0$.
 - (b) Deduce that if $m_i = 0$ for some i , then $m_1 \otimes \cdots \otimes m_n = 0$ in $M_1 \otimes_R \cdots \otimes_R M_n$.
 - (c) Prove that $2 \otimes 3 = 0$ in $\mathbb{Z}_6 \otimes_{\mathbb{Z}} \mathbb{Z}_7$.
9. Prove or disprove: if $f : M \times N \rightarrow P$ is R -linear and R -bilinear, then f must be the zero map.
10. Suppose R is any ring, possibly non-commutative. Defining R -bilinearity as in the text, prove: if M and N are left R -modules and $f : M \times N \rightarrow R$ is an R -bilinear map such that $f(x, y)$ is nonzero and not a zero divisor for some $x \in M$ and $y \in N$, then R is commutative.
11. (a) Show that if (N', i') solves the UMP in §20.5, then there exists a unique R -module isomorphism $g : \bigwedge^n M \rightarrow N'$ with $g \circ i = i'$ (where i is the map from M^n to $\bigwedge^n M$ defined in §20.5).
 - (b) State and prove a result similar to (a) for $\text{Sym}^n M$.
12. Prove: for all R -modules M , $M \cong M^{\otimes 1} \cong \bigwedge^1 M \cong \text{Sym}^1 M$.
13. Formulate and solve a UMP that converts anti-commutative maps to R -linear maps.
14. Assume $1_R + 1_R$ is not zero and not a zero divisor. (a) Let K_1 be the submodule of $M^{\otimes n}$ generated by all elements of the form

$$m_1 \otimes \cdots \otimes m_i \otimes \cdots \otimes m_j \otimes \cdots \otimes m_n + m_1 \otimes \cdots \otimes m_j \otimes \cdots \otimes m_i \otimes \cdots \otimes m_n$$

for all $m_k \in M$ and all $i < j$. Prove that $M^{\otimes n}/K_1 \cong \bigwedge^n M$. (b) Let K_2 be the submodule of $M^{\otimes n}$ generated by all elements of the form

$$m_1 \otimes \cdots \otimes m_i \otimes m_{i+1} \otimes \cdots \otimes m_n + m_1 \otimes \cdots \otimes m_{i+1} \otimes m_i \otimes \cdots \otimes m_n$$

for all $m_k \in M$ and all $i < n$. Prove that $M^{\otimes n}/K_2 \cong \bigwedge^n M$.

15. Let K_3 be the submodule of $M^{\otimes n}$ generated by all elements of the form

$$m_1 \otimes \cdots \otimes m_i \otimes m_{i+1} \otimes \cdots \otimes m_n - m_1 \otimes \cdots \otimes m_{i+1} \otimes m_i \otimes \cdots \otimes m_n$$

for all $m_k \in M$ and all $i < n$. Prove that $M^{\otimes n}/K_3 \cong \text{Sym}^n M$.

16. Assume $n! \cdot 1_R$ is invertible in R . Let L be the submodule of $M^{\otimes n}$ generated by elements of the form

$$(n!)^{-1} \sum_{w \in S_n} m_{w(1)} \otimes m_{w(2)} \otimes \cdots \otimes m_{w(n)},$$

where all $m_k \in M$. Prove $L \cong \text{Sym}^n M$.

17. Assume $n! \cdot 1_R$ is invertible in R . Let L' be the submodule of $M^{\otimes n}$ generated by elements of the form

$$(n!)^{-1} \sum_{w \in S_n} \text{sgn}(w) m_{w(1)} \otimes m_{w(2)} \otimes \cdots \otimes m_{w(n)},$$

where all $m_k \in M$. Prove $L' \cong \bigwedge^n M$.

18. Assume X_1, \dots, X_n are generating sets for the R -modules M_1, \dots, M_n (respectively). Prove that $X = \{x_1 \otimes \cdots \otimes x_n : x_i \in X_i\}$ generates the R -module $M_1 \otimes_R \cdots \otimes_R M_n$.
19. Assume X generates the R -module M . (a) Prove $\{x_1 \wedge \cdots \wedge x_n : x_i \in X\}$ generates the R -module $\bigwedge^n M$. (b) Suppose $<$ is a total ordering on X . Prove $\{x_1 \wedge \cdots \wedge x_n : x_i \in X, x_1 < \cdots < x_n\}$ generates $\bigwedge^n M$.
20. State and prove results similar to (a) and (b) in Exercise 19 for $\text{Sym}^n M$.
21. Let V be the free \mathbb{Z}_2 -module \mathbb{Z}_2^2 , which has basis $X = \{(1, 0), (0, 1)\}$. (a) Use X to find bases for $V^{\otimes 2}$, $\bigwedge^2 V$, and $\text{Sym}^2 V$. (b) List all elements $z \in V \otimes_{\mathbb{Z}_2} V$. For each z , express z in all possible ways as a basic tensor $u \otimes v$ with $u, v \in V$, or explain why this cannot be done. (c) List all $w \in \bigwedge^2 V$. For each w , express w in all possible ways in the form $u \wedge v$, or explain why this cannot be done. (d) List all $y \in \text{Sym}^2 V$. For each y , express y in all possible ways in the form uv , or explain why this cannot be done.
22. **Myth:** “The tensor product of two nonzero modules must be nonzero.” Disprove this myth by showing that for all $a, b \in \mathbb{N}^+$ with $\gcd(a, b) = 1$, $\mathbb{Z}_a \otimes_{\mathbb{Z}} \mathbb{Z}_b = \{0\}$.
23. **Myth:** “If $M \neq \{0\}$, then $M^{\otimes 2} \neq \{0\}$.” Disprove this myth by considering the \mathbb{Z} -module $M = \mathbb{Q}/\mathbb{Z}$.
24. In the proof that $M \otimes_R N \cong N \otimes_R M$ in §20.7, show directly that f is surjective.
25. *Commutativity Isomorphisms for Tensor Products.* Prove: for all R -modules M_1, \dots, M_n and all $w \in S_n$, $M_1 \otimes_R \cdots \otimes_R M_n \cong M_{w(1)} \otimes_R \cdots \otimes_R M_{w(n)}$.
26. Prove $M \otimes_R R \cong M$: (a) by using previously proved isomorphisms; (b) by defining specific isomorphisms in both directions.

27. Suppose $I \subseteq [n]$ and $M_i = R$ for all $i \in [n] \sim I$. Prove $M_1 \otimes_R \cdots \otimes_R M_n \cong \bigotimes_{i \in I} M_i$.
28. (a) Prove or disprove: for all commutative rings R , $\bigwedge^2 R \cong R$ as R -modules.
 (b) Prove or disprove: for all commutative rings R , $\text{Sym}^2 R \cong R$ as R -modules.
29. In the proof of $(M \oplus N) \otimes_R P \cong (M \otimes_R P) \oplus (N \otimes_R P)$ in §20.8: (a) check carefully that g , g_1 , and g_2 are R -bilinear; (b) check that elements of the form $(m \otimes p, 0)$ and $(0, n \otimes p)$ generate $(M \otimes_R P) \oplus (N \otimes_R P)$.
30. Prove that $P \otimes_R (M \oplus N) \cong (P \otimes_R M) \oplus (P \otimes_R N)$ by showing both sides solve the same UMP.
31. (a) Prove (20.10). (b) Prove (20.11).
32. Use the isomorphisms in §20.8 to give a new proof that the tensor product of free R -modules is free, and that (in the finite-dimensional case) $\dim(M_1 \otimes_R \cdots \otimes_R M_n) = \prod_{k=1}^n \dim(M_k)$.
33. In the proof in §20.9, check that: (a) p'_z is k -linear; (b) q'_y is m -linear; (c) t' is R -bilinear; (d) the elements $(x_1 \otimes \cdots \otimes x_k) \otimes (x_{k+1} \otimes \cdots \otimes x_{k+m})$ generate $M^{\otimes k} \otimes_R M^{\otimes m}$.
34. Carefully prove (20.12): (a) by imitating the proof in §20.9; (b) by showing both sides solve the same UMP.
35. *Extension of Scalars.* Suppose M is a free R -module with basis X , and assume R is a subring of a commutative ring S . (a) Prove $M \otimes_R S$ is a free S -module with basis $X \otimes 1_S = \{x \otimes 1_S : x \in X\}$. (b) Suppose N is another free R -module with basis Y , and $T : M \rightarrow N$ is an R -linear map with matrix A relative to the bases X and Y . Show that $T \otimes \text{id}_S : M \otimes_R S \rightarrow N \otimes_R S$ is an S -linear map with matrix A relative to the bases $X \otimes 1_S$ and $Y \otimes 1_S$.
36. (a) Given an R -linear map $f : M \rightarrow P$, give the details of the construction of the induced map $\text{Sym}^n f : \text{Sym}^n M \rightarrow \text{Sym}^n P$. (b) If also $g : P \rightarrow Q$ is R -linear, prove $(\text{Sym}^n g) \circ (\text{Sym}^n f) = \text{Sym}^n(g \circ f)$ and $\text{Sym}^n \text{id}_M = \text{id}_{\text{Sym}^n M}$.
37. The \mathbb{Z}_5 -module $V = \mathbb{Z}_5^2$ has a \mathbb{Z}_5 -basis $X = \{e_1, e_2\}$, where $e_1 = (1, 0)$ and $e_2 = (0, 1)$. Define $f : X \times X \rightarrow \mathbb{Z}_5$ by $f(e_1, e_1) = 3$, $f(e_1, e_2) = 1$, $f(e_2, e_1) = 4$, and $f(e_2, e_2) = 0$. (a) Extend f by multilinearity to $T : V \times V \rightarrow \mathbb{Z}_5$. Compute $T((2, 3), (4, 1))$. (b) T induces a \mathbb{Z}_5 -linear map $S : V \otimes_{\mathbb{Z}_5} V \rightarrow \mathbb{Z}_5$. Describe the kernel of S .
38. The \mathbb{R} -module $V = \mathbb{R}^4$ has the standard ordered basis $X = (e_1, e_2, e_3, e_4)$. (a) Write down an \mathbb{R} -basis for $V^{\otimes 2}$. (b) For all $k \geq 0$, write down an \mathbb{R} -basis for $\bigwedge^k V$. (c) Write down an \mathbb{R} -basis for $\text{Sym}^3 V$.
39. Let $X = (e_1, e_2, e_3)$ be the standard ordered basis of $V = \mathbb{R}^3$. Define $f : X^2 \rightarrow \mathbb{R}$ by $f(e_1, e_2) = 5$, $f(e_1, e_3) = -2$, and $f(e_2, e_3) = 1$. Let $T : \mathbb{R}^3 \times \mathbb{R}^3 \rightarrow \mathbb{R}$ be the unique alternating map induced from f . Compute $T((a, b, c), (x, y, z))$ for all $a, b, c, x, y, z \in \mathbb{R}$.
40. In §20.15: (a) check that T is a symmetric map; (b) prove carefully that $T = T'$; (c) show that (N', i') solves the same UMP that (N, i) does.
41. Define a bijection from the set X of weakly increasing sequences of length n using entries in $\{1, 2, \dots, d\}$ to the set Y of strictly increasing sequences of length n using entries in $\{1, 2, \dots, d+n-1\}$. Conclude that $|X| = |Y| = \binom{d+n-1}{n}$.
42. Let $A = \begin{bmatrix} 2 & 1 \\ 0 & -1 \end{bmatrix}$ and $B = \begin{bmatrix} -1 & -2 & -1 \\ 2 & 0 & 2 \\ 1 & 0 & 1 \end{bmatrix}$. Compute $A \otimes B$ and $B \otimes A$.

43. Let $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$. (a) Compute $A^{\otimes k}$ (iterated tensor product of matrices) for all $k \geq 0$. [Hint: Label the rows and columns of $A^{\otimes k}$ by subsets of $[k]$.] (b) Compute the inverses of the matrices in (a).
44. By computing the general entry on each side, prove: for $A, C \in M_m(R)$ and $B, D \in M_n(R)$, $(CA) \otimes (DB) = (C \otimes D)(A \otimes B)$.
45. (a) Fix $m, n \in \mathbb{N}^+$. Show that the map $p : M_m(R) \times M_n(R) \rightarrow M_{mn}(R)$ given by $p(A, B) = A \otimes B$ (tensor product of matrices) is R -bilinear. (b) Deduce the existence of an R -isomorphism $M_m(R) \otimes_R M_n(R) \cong M_{mn}(R)$.
46. Suppose M_1, \dots, M_n are free R -modules with respective bases X_1, \dots, X_n of sizes d_1, \dots, d_n . Suppose $f_k : M_k \rightarrow M_k$ is an R -linear map represented by the matrix $A_k \in M_{d_k}(R)$, for $1 \leq k \leq n$. Let A be the matrix of the map $\bigotimes_{k=1}^n f_k$ relative to the basis $X = \{x_1 \otimes \dots \otimes x_n : x_k \in X_k\}$ of $\bigotimes_{k=1}^n M_k$. (a) What is the size of A ? (b) Describe the entries of A in terms of the entries of the A_k 's.
47. Let $f : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ have matrix $\begin{bmatrix} 2 & 2 & -1 \\ 0 & 1 & 3 \\ 0 & -1 & 1 \end{bmatrix}$ relative to the standard ordered basis. Compute the matrix of $\bigwedge^k f$ for $0 \leq k \leq 3$.
48. Let M and N be free R -modules, and suppose an R -linear map $f : M \rightarrow N$ has matrix A relative to ordered bases X for M and Y for N . Compute the entries in the matrix of $\bigwedge^k f$ relative to the bases derived from $X_<^k$ and $Y_<^k$.
49. (a) Use Exercise 48 and the relation $\bigwedge^k(g \circ f) = (\bigwedge^k g) \circ (\bigwedge^k f)$ to deduce a formula relating the k 'th order minors of rectangular matrices A , B , and AB . (b) Deduce the Cauchy–Binet formula (§5.14) from (a).
50. In §20.18: (a) carefully check that the structure $(T(M), +, \star)$ is a ring and an R -algebra; (b) check that f_k is k -linear; (c) fill in the missing details in the proof that g is a ring homomorphism.
51. Given an R -module M , construct a commutative R -algebra $\text{Sym}(M)$ and an R -linear map $i : M \rightarrow \text{Sym}(M)$ solving this UMP: for any commutative R -algebra $(A, +, *)$ and any R -linear map $f : M \rightarrow A$, there exists a unique R -algebra homomorphism $g : \text{Sym}(M) \rightarrow A$ with $f = g \circ i$. [One can solve this problem in two ways: either imitate the construction in §20.18, or take the quotient of the tensor algebra $T(M)$ by an appropriate ideal.]
52. Given an R -module M , construct an R -algebra $\bigwedge(M)$ such that $z \star z = 0$ for all $z \in \bigwedge(M)$ and an R -linear map $i : M \rightarrow \bigwedge(M)$ solving this UMP: for any R -algebra $(A, +, *)$ and any R -linear map $f : M \rightarrow A$ such that $f(x) * f(x) = 0_A$ for all $x \in M$, there exists a unique R -algebra homomorphism $g : \bigwedge(M) \rightarrow A$ with $f = g \circ i$.
53. Suppose M , N , and P are R -modules, and $f : M \rightarrow N$, $g : N \rightarrow P$ are R -linear maps. (a) Show that f extends to a unique R -algebra homomorphism $T(f) : T(M) \rightarrow T(N)$. For $y = (y_k : k \geq 0) \in T(M)$, give a formula for $T(f)(y)$. (b) Show $T(g \circ f) = T(g) \circ T(f)$ and $T(\text{id}_M) = \text{id}_{T(M)}$. (c) Prove similar results for $\text{Sym}(M)$ and $\bigwedge(M)$ (defined in the two previous exercises).
54. *Tensor Product for Bimodules.* Let R , S , and T be rings, possibly non-commutative. We use the notation ${}_R M_S$ to indicate that M is an R, S -bimodule, which is a left R -module and a right S -module such that $(rm)s = r(ms)$ for all $r \in R$, $m \in M$, and $s \in S$. A bimodule homomorphism between two

R, S -bimodules is a map that is both R -linear and S -linear. Fix bimodules $_R M_S$, ${}_S N_T$, and ${}_R P_T$. Say that a map $g : M \times N \rightarrow P$ is S -biadditive iff $g(m + m', n) = g(m, n) + g(m', n)$, $g(m, n + n') = g(m, n) + g(m, n')$, and $g(ms, n) = g(m, sn)$ for all $m, m' \in M$, $n, n' \in N$, and $s \in S$. Call the map g an R, S, T -map iff g is S -biadditive and $g(rm, n) = rg(m, n)$ and $g(m, nt) = g(m, n)t$ for all $m \in M$, $n \in N$, $r \in R$, and $t \in T$. Construct a bimodule ${}_R(M \otimes_S N)_T$ (called the *tensor product of M and N over S*) and an R, S, T -map $i : M \times N \rightarrow M \otimes_S N$, solving the following UMP: for all ${}_R P_T$, there is a bijection from the set of group homomorphisms $f : M \otimes_S N \rightarrow P$ onto the set of S -biadditive maps $g : M \times N \rightarrow P$, that maps f to $f \circ i$. Furthermore, show that f is an R, T -bimodule homomorphism iff $g = f \circ i$ is an R, S, T -map. [Hint: Define the commutative group $M \otimes_S N$ to be a certain quotient of the free commutative group with basis $M \times N$. Verify the UMP for S -biadditive maps. Aided by this, carefully define a left R -action and a right T -action on $M \otimes_S N$, verify the bimodule axioms, and prove that bimodule homomorphisms correspond to R, S, T -maps.]

55. Establish tensor product isomorphisms (analogous to those proved in §20.8) for tensor products of bimodules.
56. Given rings R, S, T, U and bimodules ${}_R M_S$, ${}_S N_T$, ${}_T P_U$, prove there is an R, U -bimodule isomorphism $(M \otimes_S N) \otimes_T P \cong M \otimes_S (N \otimes_T P)$.
57. Let $(A, +, \star)$ be an R -algebra. (a) Show there is a well-defined R -linear map $m : A \otimes_R A \rightarrow A$ given on generators by $m(x \otimes y) = x \star y$ for $x, y \in A$. (b) Identifying $(A \otimes_R A) \otimes_R A$ and $A \otimes_R (A \otimes_R A)$ with $A^{\otimes 3}$, show that $m \circ (m \otimes \text{id}_A) = m \circ (\text{id}_A \otimes m)$. (c) Let $c : A \otimes_R A \rightarrow A \otimes_R A$ be the isomorphism given by $c(x \otimes y) = y \otimes x$ for $x, y \in A$. Show: if A is commutative, then $m \circ c = m$. (d) Define an R -linear map $e : R \rightarrow A$ by $e(r) = r1_A$ for $r \in R$. Let $g : R \otimes_R A \rightarrow A$ and $h : A \otimes_R R \rightarrow A$ be the canonical isomorphisms. Show $g = m \circ (e \otimes \text{id}_A)$ and $h = m \circ (\text{id}_A \otimes e)$.
58. Let A be an R -module. Suppose $m : A \otimes_R A \rightarrow A$ and $e : R \rightarrow A$ are R -linear maps satisfying the equations in (b) and (d) of Exercise 57. Define $x \star y = m(x \otimes y)$ for $x, y \in A$. Show that $(A, +, \star)$ is an R -algebra with identity $1_A = e(1_R)$. Also show that if m satisfies the equation in (c) of Exercise 57, then A is commutative.

This page intentionally left blank

Appendix: Basic Definitions

This appendix records some general mathematical definitions and notations that occur throughout the text. The word *iff* is defined to mean “if and only if.”

Sets

We first review some definitions from set theory that will be used constantly. All capital letters appearing below denote sets.

- *Set Membership:* $x \in S$ means x is a member of the set S .
 - *Set Non-membership:* $x \notin S$ means x is not a member of the set S .
 - *Subsets:* $A \subseteq B$ means for all x , if $x \in A$ then $x \in B$.
 - *Binary Union:* For all x , $x \in A \cup B$ iff $x \in A$ or $x \in B$.
 - *Binary Intersection:* For all x , $x \in A \cap B$ iff $x \in A$ and $x \in B$.
 - *Set Difference:* For all x , $x \in A \sim B$ iff $x \in A$ and $x \notin B$.
 - *Empty Set:* For all x , $x \notin \emptyset$.
 - *Indexed Unions:* For all x , $x \in \bigcup_{i \in I} A_i$ iff there exists $i \in I$ with $x \in A_i$.
 - *Indexed Intersections:* For all x , $x \in \bigcap_{i \in I} A_i$ iff for all $i \in I$, $x \in A_i$.
 - *Cartesian Products:* $A \times B$ is the set of all ordered pairs (a, b) with $a \in A$ and $b \in B$. More generally, $A_1 \times \cdots \times A_n$ is the set of all ordered n -tuples (a_1, \dots, a_n) with $a_i \in A_i$ for $1 \leq i \leq n$.
 - *Number Systems:* We write $\mathbb{N} = \{0, 1, 2, 3, \dots\}$ for the set of natural numbers, $\mathbb{N}^+ = \{1, 2, 3, \dots\}$ for the set of positive integers, \mathbb{Z} for the set of integers, \mathbb{Q} for the set of rational numbers, \mathbb{R} for the set of real numbers, and \mathbb{C} for the set of complex numbers. \mathbb{Q}^+ denotes the set of positive rational numbers; \mathbb{R}^+ denotes the set of positive real numbers. For all $n \in \mathbb{N}^+$, we write $[n]$ to denote the finite set $\{1, 2, \dots, n\}$.
-

Functions

Formally, a *function* is an ordered triple $f = (X, Y, G)$, where X is a set called the *domain* of f , Y is a set called the *codomain* of f , and $G \subseteq X \times Y$ is a set called the *graph* of f ,

which is required to satisfy this condition: for all $x \in X$, there exists a unique $y \in Y$ with $(x, y) \in G$. For all $x \in X$, we write $y = f(x)$ iff $(x, y) \in G$.

The notation $f : X \rightarrow Y$ means that f is a function with domain X and codomain Y . One often introduces a new function by a phrase such as: “Let $f : X \rightarrow Y$ be given by $f(x) = \dots$,” where \dots is some formula involving x . One must check that for each fixed $x \in X$, this formula always does produce exactly one output, and that this output lies in the claimed codomain Y . By our definition, two functions f and g are *equal* iff they have the same domain and the same codomain and the same graph. To check equality of the graphs, one must check that $f(x) = g(x)$ for all x in the common domain of f and g .

Given functions $f : X \rightarrow Y$ and $g : Y \rightarrow Z$, the *composite function* $g \circ f$ is the function with domain X , codomain Z , and graph $\{(x, g(f(x))) : x \in X\}$. Thus, $g \circ f : X \rightarrow Z$ satisfies $(g \circ f)(x) = g(f(x))$ for all $x \in X$.

Let $f : X \rightarrow Y$ be any function. We say f is *one-to-one* (or *injective*, or an *injection*) iff for all $x_1, x_2 \in X$, if $f(x_1) = f(x_2)$ then $x_1 = x_2$. We say f is *onto* (or *surjective*, or a *surjection*) iff for each $y \in Y$, there exists $x \in X$ with $y = f(x)$. We say f is *bijective* (or a *bijection*) iff f is one-to-one and onto iff for each $y \in Y$, there exists a unique $x \in X$ with $y = f(x)$. The composition of two injections is an injection; the composition of two surjections is a surjection; and the composition of two bijections is a bijection.

The *identity function* on any set X is the function $\text{id}_X : X \rightarrow X$ given by $\text{id}_X(x) = x$ for all $x \in X$. Given $f : X \rightarrow Y$, we say that a function $g : Y \rightarrow X$ is the *inverse* of f iff $f \circ g = \text{id}_Y$ and $g \circ f = \text{id}_X$, in which case we write $g = f^{-1}$. One can show that the inverse of f is unique when it exists; and f^{-1} exists iff f is a bijection, in which case f^{-1} is also a bijection and $(f^{-1})^{-1} = f$.

Suppose $f : X \rightarrow Y$ is a function and $Z \subseteq X$. We obtain a new function $g : Z \rightarrow Y$ with domain Z by setting $g(z) = f(z)$ for all $z \in Z$. We call g the *restriction of f to Z* , denoted $g = f|Z$ or $f|_Z$.

Suppose $f : X \rightarrow Y$ is any function. For all $A \subseteq X$, the *direct image* of A under f is the set $f[A] = \{f(a) : a \in A\} \subseteq Y$. For all $B \subseteq Y$, the *inverse image* of B under f is the set $f^{-1}[B] = \{x \in X : f(x) \in B\} \subseteq X$. This notation is not meant to suggest that the inverse function f^{-1} must exist. But, when f^{-1} does exist, the inverse image of B under f coincides with the direct image of B under f^{-1} , so the notation $f^{-1}[B]$ is not ambiguous. The *image* of f is the set $f[X]$; f is a surjection iff $f[X] = Y$. We use square brackets for direct and inverse images to prevent ambiguity. More precisely, if A is both a member of X and a subset of X , then $f(A)$ is the value of f at the point A in its domain, whereas $f[A]$ is the direct image under f of the subset A of the domain.

Relations

A *relation from X to Y* is a subset R of $X \times Y$. For $x \in X$ and $y \in Y$, xRy means $(x, y) \in R$. A *relation on a set X* is a relation R from X to X . R is called *reflexive on X* iff for all $x \in X$, xRx . R is called *symmetric* iff for all x, y , xRy implies yRx . R is called *antisymmetric* iff for all x, y , xRy and yRx implies $x = y$. R is called *transitive* iff for all x, y, z , xRy and yRz implies xRz . R is called an *equivalence relation on X* iff R is reflexive on X , symmetric, and transitive.

Suppose R is an equivalence relation on a set X . For $x \in X$, the *equivalence class of x relative to R* is the set $[x]_R = \{y \in X : xRy\}$. A given equivalence class typically has many names; more precisely, for all $x, z \in R$, $[x]_R = [z]_R$ iff xRz . The *quotient set X modulo R* is the set of all equivalence classes of R , namely $X/R = \{[x]_R : x \in X\}$.

A *set partition* of a given set X is a collection P of nonempty subsets of X such that for all $x \in X$, there exists a unique $S \in P$ with $x \in S$. For every equivalence relation R on a fixed set X , the quotient set X/R is a set partition of X consisting of the equivalence classes $[x]_R$ for $x \in X$. Conversely, given any set partition P of X , the relation R defined by “ xRy iff there exists $S \in P$ with $x \in S$ and $y \in S$ ” is an equivalence relation on X with $X/R = P$. Formally, letting EQ_X be the set of all equivalence relations on X and letting SP_X be the set of all set partitions on X , the map $f : EQ_X \rightarrow SP_X$ given by $f(R) = X/R$ for $R \in EQ_X$ is a bijection.

Partially Ordered Sets

A *partial ordering* on a set X is a relation \leq on X that is reflexive on X , antisymmetric, and transitive. A *partially ordered set* or *poset* is a pair (X, \leq) where \leq is a partial ordering on X . A poset (X, \leq) is *totally ordered* iff for all $x, y \in X$, $x \leq y$ or $y \leq x$. More generally, a subset Y of a poset (X, \leq) is called a *chain* iff for all $x, y \in Y$, $x \leq y$ or $y \leq x$. By definition, $y \geq x$ means $x \leq y$; $x < y$ means $x \leq y$ and $x \neq y$; and $x > y$ means $x \geq y$ and $x \neq y$.

Let S be a subset of a poset (X, \leq) . An *upper bound* for S is an element $x \in X$ such that for all $y \in S$, $y \leq x$. A *greatest element* of S is an element $x \in S$ such that for all $y \in S$, $y \leq x$; x is unique if it exists. A *lower bound* for S is an element $x \in X$ such that for all $y \in S$, $x \leq y$. A *least element* of S is an element $x \in S$ such that for all $y \in S$, $x \leq y$; x is unique if it exists.

We say $x \in X$ is a *least upper bound* for S iff x is the least element of the set of upper bounds of S in X . In detail, this means $y \leq x$ for all $y \in S$; and for any $z \in X$ such that $y \leq z$ for all $y \in S$, $x \leq z$. The least upper bound of S is unique if it exists; we write $x = \sup S$ in this case. For $S = \{y_1, \dots, y_n\}$, the notation $y_1 \vee y_2 \vee \dots \vee y_n$ is also used to denote $\sup S$.

We say $x \in X$ is a *greatest lower bound* for S iff x is the greatest element of the set of lower bounds of S in X . In detail, this means $x \leq y$ for all $y \in S$; and for any $z \in X$ such that $z \leq y$ for all $y \in S$, $z \leq x$. The greatest lower bound of S is unique if it exists; we write $x = \inf S$ in this case. For $S = \{y_1, \dots, y_n\}$, the notation $y_1 \wedge y_2 \wedge \dots \wedge y_n$ is also used to denote $\inf S$.

A *lattice* is a poset (X, \leq) such that for all $a, b \in X$, the least upper bound $a \vee b$ and the greatest lower bound $a \wedge b$ exist in X . A *complete lattice* is a poset (X, \leq) such that for every nonempty subset S of X , $\sup S$ and $\inf S$ exist in X .

A *maximal element* in a poset (X, \leq) is an element $x \in X$ such that for all $y \in X$, if $x \leq y$ then $y = x$. A *minimal element* of X is an element $x \in X$ such that for all $y \in X$, if $y \leq x$ then $y = x$. Zorn's lemma states that if (X, \leq) is a poset in which every chain $Y \subseteq X$ has an upper bound in X , then X has a maximal element. Zorn's lemma is discussed in detail in §16.6.

This page intentionally left blank

Further Reading

Chapter 1. There are many introductory accounts of modern algebra, including the texts by Durbin [14], Fraleigh [15], Gallian [16], and Rotman [48]. For more advanced treatments of modern algebra, one may consult the textbooks by Dummit and Foote [13], Hungerford [29], Jacobson [30], and Rotman [47]. Introductions to linear algebra at various levels abound; among many others, we mention the books by Larson and Falvo [33], Lay [34], and Strang [56]. Two more advanced linear algebra books that are similar, in some respects, to the present volume are the texts by Halmos [24] and Hoffman and Kunze [27].

Chapter 2. For basic facts on permutations, one may consult any of the abstract algebra texts mentioned above. There is a vast literature on permutations and the symmetric group; we direct the reader to the texts by Bona [6], Rotman [51], and Sagan [54] for more information.

Chapter 3. Thorough algebraic treatments of polynomials may be found in most texts on abstract algebra, such as those by Dummit and Foote [13] or Hungerford [29]. For more details on formal power series, one may consult [36, Chpt. 7]. The matrix reduction algorithm in §3.10 for computing gcds of polynomials (or integers) comes from an article by W. Blankinship [5]. Cox, Little, and O’Shea have written an excellent book [11] on multivariable polynomials and their role in computational algebraic geometry.

Chapter 4. Three classic texts on matrix theory are the books by Gantmacher [18], Horn and Johnson [28], and Lancaster [32].

Chapter 5. There is a vast mathematical literature on the subject of determinants. Lacking the space to give tribute to all of these, we only mention the text by Turnbull [59], the book of Aitken [2], and the treatise of Muir [41]. Muir has also written an extensive four-volume work chronicling the historical development of the theory of determinants [40].

Chapter 6. This chapter developed a “dictionary” linking abstract concepts defined for vector spaces and linear maps to concrete concepts defined for column vectors and matrices. Many texts on matrix theory, such as Horn and Johnson [28], heavily favor matrix-based descriptions and proofs. Other texts, most notably Bourbaki [7, Chpt. II], prefer a very abstract development that makes almost no mention of matrices. We think it is advisable to gain facility with both languages for discussing linear algebra. For alternative developments of this material, see Halmos [24] and Hoffman and Kunze [27].

Chapter 7. The analogy between complex numbers and complex matrices, including the theorems on the polar decomposition of a matrix, is based on the exposition in Halmos [24]. A wealth of additional material on properties of Hermitian, unitary, positive definite, and normal matrices may be found in Horn and Johnson’s text [28].

Chapter 8. One can derive the Jordan canonical form theorem in many different ways. In abstract algebra, one often deduces this theorem from the rational canonical form

theorem [27], which in turn is derivable from the classification of finitely generated modules over principal ideal domains. See Chapter 18 or [30] for this approach. Matrix theorists might prefer a more algorithmic construction that triangularizes a complex matrix and then gradually reduces it to Jordan form [12, 28]. Various elementary derivations can be found in [8, 17, 19, 23, 60].

Chapter 9. Further information on QR factorizations can be found in Chapter 5 of Golub and van Loan [20], part II of Trefethen and Bau [58], and §5.3 of Kincaid and Cheney [31]. For LU factorizations, see [20, Chpt. 3] or [28, Sec. 3.5]. These references also contain a wealth of information on the numerical stability properties of matrix factorizations and the associated algorithms.

Chapter 10. Our treatment of iterative algorithms for solving linear systems and computing eigenvalues is similar to that found in §4.6 and §5.1 of Kincaid and Cheney [31]. For more information on this topic, one may consult the numerical analysis texts authored by Ackleh, Allen, Kearfott, and Seshaiyer [1, §3.4, Chpt. 5], Cheney and Kincaid [9, §8.2, §8.4], Trefethen and Bau [58, Chpt. VI], and Golub and Van Loan [20].

Chapter 11. For a very detailed treatment of convex sets and convex functions, the reader may consult Rockafellar's text [46]. A wealth of material on convex polytopes can be found in Grünbaum's encyclopedic work [21]. The presentation in §11.17 through §11.21 is similar to [63, Lecture 1].

Chapter 12. Our exposition of ruler and compass constructions is similar to the accounts found in [30, Vol. 1, Chpt. 4] and [49, App. C]. Treatments of Galois theory may be found in these two texts, as well as in the books by Cox [10], Dummit and Foote [13], and Hungerford [29]. See Tignol's book [57] for a very nice historical account of the development of Galois theory. Another good reference for geometric constructions and other problems in field theory is Hadlock [22].

Chapter 13. Another treatment of dual spaces and their relation to complex inner product spaces appears in Halmos [24]. A good discussion of dual spaces in the context of Banach spaces is given in Simmons [55, Chpt. 9]. The book of Cox, Little, and O'Shea [11] contains an excellent exposition of the ideal-variety correspondence and other aspects of affine algebraic geometry.

Chapter 14. Two other introductions to Hilbert spaces at a level similar to ours can be found in Simmons [55, Chpt. 10] and Rudin [53, Chpt. 4]. Halmos' book [25] contains an abundance of problems on Hilbert spaces. For more on metric spaces, see Simmons [55] or Munkres [42].

Chapter 15. A nice treatment of commutative groups, finitely generated or not, appears in Chapter 10 of Rotman's group theory text [51]. Another exposition of the reduction algorithm for integer matrices and its connection to classifying finitely generated commutative groups is given by Munkres [43, §11].

Chapter 16. Our treatment of independence structures is based on [30, Vol. 2]. A very nice introduction to matroids is [62, Sec. 8.2]. The books by Oxley [45] and Welsh [61] contain detailed accounts of matroid theory.

Chapter 17. Four excellent accounts of module theory appear in Anderson and Fuller's

text [3], Atiyah and Macdonald's book [4], Jacobson's *Basic Algebra 1 and 2* [30] (especially the third chapter in each volume), and Rotman's homological algebra book [52]. Bourbaki [7, Chpt. II] provides a very thorough and general, but rather difficult, treatment of modules.

Chapter 18. The classification of finitely generated modules over principal ideal domains is a standard topic covered in advanced abstract algebra texts such as [13, 30]. We hope that our coverage may be more quickly accessible to readers with a little less background in group theory and ring theory. There are two approaches to proving the rational canonical form for square matrices over a field. The approach adopted here deduces this result from the general theory for PIDs. The other approach avoids the abstraction of PIDs by proving all necessary results at the level of finite-dimensional vector spaces, T -invariant subspaces, and T -cyclic subspaces. See [27] for such a treatment. The author's opinion is that proving the special case of the classification theorem for torsion $F[x]$ -modules is not much simpler than proving the full theorem for all finitely generated modules over all PIDs. In fact, because of all the extra structure of the ring $F[x]$, focusing on this special case might even give the reader less intuition for what the proof is doing. To help the reader build intuition, we chose to cover the much more concrete case of \mathbb{Z} -modules in an earlier chapter.

Chapter 19. Two sources that give due emphasis to the central role of universal mapping properties in abstract algebra are Jacobson's two-volume algebra text [30] and Rotman's homological algebra book [52]. The appropriate general context for understanding UMP's is category theory, the basic elements of which are covered in the two references just cited. A more comprehensive introduction to category theory is given in Mac Lane's book [38].

Chapter 20. A nice introduction to multilinear algebra is the text by Northcott [44]. A very thorough account of the subject, including detailed discussions of tensor algebras, exterior algebras, and symmetric algebras, appears in [7, Chpt. III].

This page intentionally left blank

Bibliography

- [1] Azmy Ackleh, Edward J. Allen, Ralph Kearfott, and Padmanabhan Seshaiyer, *Classical and Modern Numerical Analysis: Theory, Methods, and Practice*, Chapman and Hall/CRC Press, Boca Raton, FL (2010).
- [2] A. C. Aitken, *Determinants and Matrices* (eighth ed.), Oliver and Boyd Ltd., Edinburgh (1954).
- [3] Frank W. Anderson and Kent R. Fuller, *Rings and Categories of Modules* (Graduate Texts in Mathematics Vol. 13, second ed.), Springer-Verlag, New York (1992).
- [4] M. F. Atiyah and I. G. Macdonald, *Introduction to Commutative Algebra*, Addison-Wesley, Reading, MA (1969).
- [5] W. A. Blankinship, “A new version of the Euclidean algorithm,” *Amer. Math. Monthly* **70** #7 (1963), 742–745.
- [6] Miklós Bóna, *Combinatorics of Permutations*, Chapman and Hall/CRC, Boca Raton, FL (2004).
- [7] Nicolas Bourbaki, *Algebra 1*, Springer-Verlag, New York (1989).
- [8] R. Brualdi, “The Jordan canonical form: an old proof,” *Amer. Math. Monthly* **94** #3 (1987), 257–267.
- [9] E. Ward Cheney and David R. Kincaid, *Numerical Mathematics and Computing* (sixth ed.), Brooks/Cole, Pacific Grove, CA (2007).
- [10] David A. Cox, *Galois Theory* (second ed.), John Wiley and Sons, New York (2012).
- [11] David A. Cox, John Little, and Donal O’Shea, *Ideals, Varieties, and Algorithms: An Introduction to Computational Algebraic Geometry and Commutative Algebra* (third ed.), Springer-Verlag, New York (2010).
- [12] R. Fletcher and D. Sorenson, “An algorithmic derivation of the Jordan canonical form,” *Amer. Math. Monthly* **90** #1 (1983), 12–16.
- [13] David S. Dummit and Richard M. Foote, *Abstract Algebra* (third ed.), John Wiley and Sons, New York (2003).
- [14] John R. Durbin, *Modern Algebra: An Introduction* (sixth ed.), John Wiley and Sons, New York (2008).
- [15] John B. Fraleigh, *A First Course in Abstract Algebra* (seventh ed.), Addison Wesley, Reading (2002).
- [16] Joseph A. Gallian, *Contemporary Abstract Algebra* (fifth ed.), Houghton Mifflin, Boston (2001).

- [17] A. Galperin and Z. Waksman, “An elementary approach to Jordan theory,” *Amer. Math. Monthly* **87** #9 (1980), 728–732.
- [18] F. R. Gantmacher, *The Theory of Matrices* (two volumes), Chelsea Publishing Co., New York (1960).
- [19] I. Gohberg and S. Goldberg, “A simple proof of the Jordan decomposition theorem for matrices,” *Amer. Math. Monthly* **103** #2 (1996), 157–159.
- [20] Gene Golub and Charles Van Loan, *Matrix Computations* (third ed.), The Johns Hopkins University Press, Baltimore (1996).
- [21] Branko Grünbaum, *Convex Polytopes* (Graduate Texts in Mathematics Vol. 221, second ed.), Springer-Verlag, New York (2003).
- [22] Charles R. Hadlock, *Field Theory and Its Classical Problems*, Carus Mathematical Monograph no. 19, Mathematical Association of America, Washington, D.C., (1978).
- [23] J. Hall, “Another elementary approach to the Jordan form,” *Amer. Math. Monthly* **98** #4 (1991), 336–340.
- [24] Paul R. Halmos, *Finite-Dimensional Vector Spaces*, Springer-Verlag, New York (1974).
- [25] Paul R. Halmos, *A Hilbert Space Problem Book* (Graduate Texts in Mathematics Vol. 19, second ed.), Springer-Verlag, New York (1982).
- [26] Paul R. Halmos, *Naive Set Theory*, Springer-Verlag, New York (1998).
- [27] Kenneth Hoffman and Ray Kunze, *Linear Algebra* (second ed.), Prentice Hall, Upper Saddle River, NJ (1971).
- [28] Roger Horn and Charles Johnson, *Matrix Analysis* (second ed.), Cambridge University Press, Cambridge (2012).
- [29] Thomas W. Hungerford, *Algebra* (Graduate Texts in Mathematics Vol. 73), Springer-Verlag, New York (1980).
- [30] Nathan Jacobson, *Basic Algebra I and II* (second ed.), Dover Publications, Mineola, NY (2009).
- [31] David Kincaid and Ward Cheney, *Numerical Analysis: Mathematics of Scientific Computing* (second ed.), Brooks/Cole, Pacific Grove, CA (1996).
- [32] Peter Lancaster, *Theory of Matrices*, Academic Press, New York (1969).
- [33] Ron Larson and David Falvo, *Elementary Linear Algebra* (sixth ed.), Brooks Cole, Belmont, CA (2009).
- [34] David C. Lay, *Linear Algebra and Its Applications* (fourth ed.), Addison Wesley, Reading, MA (2011).
- [35] Hans Liebeck, “A proof of the equality of column and row rank of a matrix,” *Amer. Math. Monthly* **73** #10 (1966), 1114.
- [36] Nicholas A. Loehr, *Bijective Combinatorics*, Chapman and Hall/CRC, Boca Raton, FL (2011).

- [37] Nicholas A. Loehr, "A direct proof that row rank equals column rank," *College Math. J.* **38** #4 (2007), 300–301.
- [38] Saunders Mac Lane, *Categories for the Working Mathematician* (Graduate Texts in Mathematics Vol. 5, second ed.), Springer-Verlag, New York (1998).
- [39] J. Donald Monk, *Introduction to Set Theory*, McGraw-Hill, New York (1969).
- [40] Thomas Muir, *The Theory of Determinants in the Historical Order of Development* (four volumes), Dover Publications, New York (1960).
- [41] Thomas Muir, *A Treatise on the Theory of Determinants*, revised and enlarged by William Metzler, Dover Publications, New York (1960).
- [42] James R. Munkres, *Topology* (second ed.), Prentice Hall, Upper Saddle River, NJ (2000).
- [43] James R. Munkres, *Elements of Algebraic Topology*, Perseus Publishing, Cambridge, MA (1984).
- [44] D. G. Northcott, *Multilinear Algebra*, Cambridge University Press, Cambridge (1984).
- [45] James G. Oxley, *Matroid Theory* (second ed.), Oxford University Press, Oxford (2011).
- [46] R. Tyrrell Rockafellar, *Convex Analysis*, Princeton University Press, Princeton, NJ (1972).
- [47] Joseph J. Rotman, *Advanced Modern Algebra* (second ed.), American Mathematical Society, Providence, RI (2010).
- [48] Joseph J. Rotman, *A First Course in Abstract Algebra* (third ed.), Prentice Hall, Upper Saddle River, NJ (2005).
- [49] Joseph J. Rotman, *Galois Theory* (second ed.), Springer-Verlag, New York (1998).
- [50] Joseph J. Rotman, *An Introduction to Algebraic Topology* (Graduate Texts in Mathematics, Vol. 119), Springer-Verlag, New York (1988).
- [51] Joseph J. Rotman, *An Introduction to the Theory of Groups* (fourth ed.), Springer-Verlag, New York (1994).
- [52] Joseph J. Rotman, *Notes on Homological Algebra*, Van Nostrand Reinhold, New York (1970).
- [53] Walter Rudin, *Real and Complex Analysis* (third ed.), McGraw-Hill, Boston (1987).
- [54] Bruce E. Sagan, *The Symmetric Group: Representations, Combinatorial Algorithms, and Symmetric Functions* (second ed.), Springer-Verlag, New York (2001).
- [55] George F. Simmons, *Introduction to Topology and Modern Analysis*, Krieger Publishing Co., Malabar, FL (2003).
- [56] Gilbert Strang, *Introduction to Linear Algebra* (fourth ed.), Wellesley Cambridge Press, Wellesley, MA (2009).
- [57] Jean-Pierre Tignol, *Galois' Theory of Algebraic Equations*, World Scientific Publishing, Singapore (2001).

- [58] Lloyd N. Trefethen and David Bau III, *Numerical Linear Algebra*, SIAM, Philadelphia (1997).
- [59] Herbert W. Turnbull, *The Theory of Determinants, Matrices, and Invariants* (second ed.), Blackie and Son Ltd., London (1945).
- [60] H. Valiaho, “An elementary approach to the Jordan form of a matrix,” *Amer. Math. Monthly* **93** #9 (1986), 711–714.
- [61] D. J. Welsh, *Matroid Theory*, Academic Press, New York (1976).
- [62] Douglas B. West, *Introduction to Graph Theory* (second ed.), Prentice Hall, Upper Saddle River, NJ (2001).
- [63] Günter M. Ziegler, *Lectures on Polytopes* (Graduate Texts in Mathematics, Vol. 152), Springer-Verlag, New York (1995).

This page intentionally left blank

TEXTBOOKS in MATHEMATICS

Advanced Linear Algebra explores a variety of advanced topics in linear algebra that highlight the rich interconnections of the subject to geometry, algebra, analysis, combinatorics, numerical computation, and many other areas of mathematics. The book's 20 chapters are grouped into six main subjects: algebraic structures, matrices, structured matrices, geometric aspects of linear algebra, modules, and multilinear algebra. The level of abstraction gradually increases as you proceed through the text, moving from matrices to vector spaces to modules.

Each chapter consists of a mathematical vignette devoted to the development of one specific topic. Some chapters look at introductory material from a sophisticated or abstract viewpoint while others provide elementary expositions of more theoretical concepts. Several chapters offer unusual perspectives or novel treatments of standard results. Unlike similar advanced mathematical texts, this one minimizes the dependence of each chapter on material found in previous chapters so that you may immediately turn to the relevant chapter without first wading through pages of earlier material to access the necessary algebraic background and theorems.

Features

- Focuses on theoretical aspects of linear algebra
- Gives complete proofs of all results
- Supplements the general theory with many specific examples and concrete computations
- Requires very little knowledge of abstract algebra beyond fields, vector spaces over a field, subspaces, linear transformations, linear independence, and bases
- Includes a summary and exercise sets at the end of each chapter



CRC Press
Taylor & Francis Group
an informa business
www.crcpress.com

6000 Broken Sound Parkway, NW
Suite 300, Boca Raton, FL 33487
711 Third Avenue
New York, NY 10017
2 Park Square, Milton Park
Abingdon, Oxon OX14 4RN, UK

K15529

ISBN : 978-1-4665-5901-1

90000



9 781466 559011