

Introduction to Modern Algebra

David Joyce
Clark University

Version 1.2.7, 5 Dec 2017 ¹

¹Copyright (C) 2008,2017.

I dedicate this book to my friend and colleague Arthur Chou. Arthur encouraged me to write this book. I'm sorry that he did not live to see it finished.

Arthur was born in 1954 in Taipei, Taiwan. He received his bachelors in mathematics in 1976 from Tunghai University and his PhD from Stony Brook in 1982. After a year at the Institute for Advanced Study at Princeton, he joined Clark University in 1983. He was promoted to Associate Professor six years later and promoted to full professor in 2008, the year he died. Besides mathematics, he had many other interests. Among other things, he was the general manager of the North America Elite Youth Orchestra which performed at Dallas, Beijing, and Taipei, and he was the deacon of the Chinese Gospel Church in Southborough, Massachusetts.

Contents

List of Figures	vi
List of Tables	viii
1 Introduction	1
1.1 Algebra	1
1.2 Structures in Modern Algebra	2
1.2.1 Operations on sets	2
1.2.2 Fields	4
1.2.3 Rings	5
1.2.4 Groups	6
1.2.5 Other algebraic structures besides fields, rings, and groups	10
1.3 Isomorphisms, homomorphisms, etc.	11
1.3.1 Isomorphisms	11
1.3.2 Homomorphisms	12
1.3.3 Monomorphisms and epimorphisms	13
1.3.4 Endomorphisms and automorphisms	14
1.4 A little number theory	15
1.4.1 Mathematical induction on the natural numbers \mathbf{N}	15
1.4.2 Divisibility	16
1.4.3 Prime numbers	17
1.4.4 The Euclidean algorithm	19
1.5 The fundamental theorem of arithmetic	22
1.6 Polynomials	25
1.6.1 Division for polynomials	26
1.6.2 Roots of unity and cyclotomic polynomials	28
2 Fields	31
2.1 Introduction to fields	31
2.1.1 Definition of fields	31
2.1.2 Subtraction, division, multiples, and powers	32
2.1.3 Properties that follow from the axioms	33
2.1.4 Subfields	34
2.1.5 Fields of rational functions	35
2.1.6 Vector spaces over arbitrary fields	35
2.2 Cyclic rings and finite fields	35

2.2.1	The cyclic ring \mathbf{Z}_n	36
2.2.2	The cyclic prime fields \mathbf{Z}_p	39
2.2.3	Characteristics of fields, and prime fields	41
2.3	Field Extensions, algebraic fields, the complex numbers	41
2.3.1	Algebraic fields	42
2.3.2	The field of complex numbers \mathbf{C}	43
2.3.3	General quadratic extensions	44
2.4	Real numbers and ordered fields	45
2.4.1	Ordered fields	45
2.4.2	Archimedean orders	47
2.4.3	Complete ordered fields	49
2.5	Skew fields (division rings) and the quaternions	50
2.5.1	Skew fields (division rings)	50
2.5.2	The quaternions \mathbf{H}	51
3	Rings	55
3.1	Introduction to rings	55
3.1.1	Definition and properties of rings	55
3.1.2	Products of rings	57
3.1.3	Integral domains	57
3.1.4	The Gaussian integers, $\mathbf{Z}[i]$	59
3.1.5	Finite fields again	59
3.2	Factoring \mathbf{Z}_n by the Chinese remainder theorem	60
3.2.1	The Chinese remainder theorem	60
3.2.2	Brahmagupta's solution	62
3.2.3	Qin Jiushao's solution	62
3.3	Boolean rings	63
3.3.1	Introduction to Boolean rings	64
3.3.2	Factoring Boolean rings	65
3.3.3	A partial order on a Boolean ring	66
3.4	The field of rational numbers, fields of fractions	67
3.5	Categories and the category of rings	69
3.5.1	The formal definition of categories	70
3.5.2	The category \mathcal{R} of rings	71
3.5.3	Monomorphisms and epimorphisms in a category	73
3.6	Kernels, ideals, and quotient rings	74
3.6.1	Kernels of ring homomorphisms	74
3.6.2	Ideals of a ring	74
3.6.3	Quotient rings, R/I	76
3.6.4	Prime and maximal ideals	78
3.7	Krull's theorem	79
3.8	UFDs, PIDs, and EDs	80
3.8.1	Divisibility in an integral domain	80
3.8.2	Unique factorization domains	81
3.8.3	Principal ideal domains	82
3.8.4	Euclidean domains	84

3.9	Real and complex polynomial rings $\mathbf{R}[x]$ and $\mathbf{C}[x]$	87
3.9.1	$\mathbf{C}[x]$ and the Fundamental Theorem of Algebra	87
3.9.2	The polynomial ring $\mathbf{R}[x]$	89
3.10	Rational and integer polynomial rings	90
3.10.1	Roots of polynomials	90
3.10.2	Gauss's lemma and Eisenstein's criterion	92
3.10.3	Prime cyclotomic polynomials	95
3.10.4	Polynomial rings with coefficients in a UFD, and polynomial rings in several variables.	95
3.11	Number fields and their rings of integers	97
4	Groups	99
4.1	Groups and subgroups	99
4.1.1	Definition and basic properties of groups	99
4.1.2	Subgroups	100
4.1.3	Cyclic groups and subgroups	101
4.1.4	Products of groups	102
4.1.5	Cosets and Lagrange's theorem	103
4.2	Symmetric Groups S_n	104
4.2.1	Permutations and the symmetric group	104
4.2.2	Even and odd permutations	106
4.2.3	Alternating and dihedral groups	107
4.3	Cayley's theorem and Cayley graphs	110
4.3.1	Cayley's theorem	110
4.3.2	Some small finite groups	112
4.4	The category of groups \mathcal{G}	115
4.5	Conjugacy classes and quandles	115
4.5.1	Conjugacy classes	116
4.5.2	Quandles and the operation of conjugation	117
4.6	Kernels, normal subgroups, and quotient groups	120
4.6.1	Kernels of group homomorphisms and normal subgroups	120
4.6.2	Quotient groups, and projections $\gamma : G \rightarrow G/N$	121
4.6.3	Isomorphism theorems	122
4.6.4	Internal direct products	123
4.7	Matrix rings and linear groups	124
4.7.1	Linear transformations	124
4.7.2	The general linear groups $GL_n(R)$	125
4.7.3	Other linear groups	126
4.7.4	Projective space and the projective linear group $PGL_n(F)$	127
4.8	Structure of finite groups	130
4.8.1	Simple groups	131
4.8.2	The Jordan-Hölder theorem	131
4.9	Abelian groups	134
4.9.1	The category \mathcal{A} of Abelian groups	136
4.9.2	Finite Abelian groups	137

Appendices	141
A Background mathematics	143
A.1 Logic and proofs	143
A.2 Sets	144
A.2.1 Basic set theory	144
A.2.2 Functions and relations	149
A.2.3 Equivalence relations	150
A.2.4 Axioms of set theory	151
A.3 Ordered structures	153
A.3.1 Partial orders and posets.	153
A.3.2 Lattices	155
A.3.3 Boolean algebras.	157
A.4 Axiom of choice	158
A.4.1 Zorn's lemma	158
A.4.2 Well-ordering principle	159
Index	161

List of Figures

1.1	Equilateral triangle with lines of symmetry	8
1.2	Unit circle S^1	10
1.3	Divisors of 432	17
1.4	Divisibility up through 12	17
2.1	Cyclic rings $\mathbf{Z}_6, \mathbf{Z}_{19}, \mathbf{Z}$	37
3.1	Lattice of Gaussian integers $\mathbf{Z}[i]$	59
3.2	Free Boolean ring on two elements	66
3.3	Lattice of Eisenstein integers	85
3.4	Primitive 7 th roots of unity	95
4.1	Subgroups of S_3	107
4.2	Symmetries of a pentagon	108
4.3	Symmetries of a cube and tetrahedron	109
4.4	Cayley graph for D_5	111
4.5	Cayley graph for A_4	113
4.6	Distributivity in a involutory quandle	118
4.7	A conjugacy class in the quaternion group	119
4.8	The conjugacy class of transpositions in S_4	120
4.9	The Fano plane \mathbf{Z}_2P^2	128
4.10	The projective plane \mathbf{Z}_3P^2	129
4.11	Cayley graph of the Frobenius group $F_{21} = C_7 \rtimes C_3$	134
4.12	Heptahedron on a torus	135
A.1	Lattice of the Powerset of 4 elements	157

List of Tables

1.1	Composition table for six particular rational functions.	9
3.1	Notations in Boolean algebras, set theory, and Boolean rings.	65
4.1	List of small groups	114
A.1	Standard logical symbols	144

Chapter 1

Introduction

1.1 Algebra

The word “algebra” means many things. The word dates back about 1200 years ago to part of the title of al-Khwārizmī’s book on the subject, but the subject itself goes back 4000 years ago to ancient Babylonia and Egypt. It was about solving numerical problems that we would now identify as linear and quadratic equations. Versions of the quadratic formula were used to find solutions to those quadratic equations. Al-Khwārizmī (ca. 780–ca. 850) codified the algorithms (“algorithm” is a word derived from his name) for solving these equations. He wrote all his equations out in words since symbolic algebra had yet to be invented.

Other places in the world also had algebra and developed various aspects of it. The ancient Chinese solved systems of simultaneous linear equations and later developed algorithms to find roots of polynomials of high degree. Various aspects of number theory were studied in China, in India, and by Greek mathematicians.

Symbolic algebra was developed in the 1500s. Symbolic algebra has symbols for the arithmetic operations of addition, subtraction, multiplication, division, powers, and roots as well as symbols for grouping expressions (such as parentheses), and most importantly, used letters for variables.

Once symbolic algebra was developed in the 1500s, mathematics flourished in the 1600s. Coordinates, analytic geometry, and calculus with derivatives, integrals, and series were developed in that century.

Algebra became more general and more abstract in the 1800s as more algebraic structures were invented. Hamilton (1805–1865) invented quaternions (see section 2.5.2) and Grassmann (1809–1977) developed exterior algebras in the 1840s, both of which led to vector spaces. (See section 2.1.6 for vector spaces.)

Groups were developed over the 1800s, first as particular groups of substitutions or permutations, then in the 1850’s Cayley (1821–1895) gave the general definition for a group. (See chapter 2 for groups.)

Several fields were studied in mathematics for some time including the field of real numbers the field of rational number, and the field of complex numbers, but there was no general definition for a field until the late 1800s. (See chapter 2 for fields.)

Rings also were studied in the 1800s. Noether (1882–1935) gave general concept of commutative ring in 1921 which was later generalized to include noncommutative rings. (See

chapter 3 for rings.)

We'll introduce the concepts of field, ring, and group in the Introduction, then study each in turn in the following chapters.

1.2 Structures in Modern Algebra

Fields, rings, and groups. We'll be looking at several kinds of algebraic structures this semester, the three major kinds being fields in chapter 2, rings in chapter 3, and groups in chapter 4, but also minor variants of these structures.

We'll start by examining the definitions and looking at some examples. For the time being, we won't prove anything; that will come in later chapters when we look at those structures in depth.

A note on notation. We'll use the standard notation for various kinds of numbers. The set of natural numbers, $\{0, 1, 2, \dots\}$ is denoted \mathbf{N} . The set of integers $\{\dots, -2, -1, 0, 1, 2, \dots\}$ is denoted \mathbf{Z} (for *Zahlen*, German for whole number). The set of rational numbers, that is, numbers of the form $\frac{m}{n}$ where m is an integer and n is a nonzero integer, is denoted \mathbf{Q} (for "quotient"). The set of all real numbers, including all positive numbers, all negative numbers, and 0, is denoted \mathbf{R} . And the set of complex numbers, that is, numbers of the form $x + iy$ where x and y are real numbers and $i^2 = -1$, is denoted \mathbf{C} .

1.2.1 Operations on sets

For background on sets, see the section A.2 in the appendix.

We're familiar with many operations on the real numbers \mathbf{R} —addition, subtraction, multiplication, division, negation, reciprocation, powers, roots, etc.

Addition, subtraction, and multiplication are examples of binary operations, that is, functions $\mathbf{R} \times \mathbf{R} \rightarrow \mathbf{R}$ which take two real numbers as their arguments and return another real number. Division is almost a binary operation, but since division by 0 is not defined, it's only a partially defined binary operation. Most of our operations will be defined everywhere, but some, like division, won't be.

Negation is a unary operation, that is, a function $\mathbf{R} \rightarrow \mathbf{R}$ which takes one real number as an argument and returns a real number. Reciprocation is a partial unary operation since the reciprocal of zero is not defined.

The operations we'll consider are all binary or unary. Ternary operations can certainly be defined, but useful ternary operations are rare.

Some of these operations satisfy familiar identities. For example, addition and multiplication are both commutative; they satisfy the identities

$$x + y = y + x \quad \text{and} \quad xy = yx.$$

A binary operation is said to be *commutative* when the order that the two arguments are applied doesn't matter, that is, interchanging them, or commuting one across the other, doesn't change the result. Subtraction and division, however, are not commutative.

Addition and multiplication are also associative binary operations

$$(x + y) + z = x + (y + z) \quad \text{and} \quad (xy)z = x(yz).$$

A binary operation is said to be *associative* when the parentheses can be associated with either the first pair or the second pair when the operation is applied to three arguments and the result is the same. Neither subtraction nor division are associative.

Both addition and multiplication also have identity elements

$$0 + x = x = x + 0 \quad \text{and} \quad 1x = x = x1.$$

An *identity element*, also called a *neutral element*, for a binary operation is an element in the set that doesn't change the value of other elements when combined with them under the operation. So, 0 is the identity element for addition, and 1 is the identity element for multiplication. Subtraction and division don't have identity elements. (Well, they do on the right, since $x - 0 = x$ and $\frac{x}{1} = x$, but not on the left, since usually $0 - x \neq x$ and $\frac{1}{x} \neq x$.)

Also, there are additive inverses and multiplicative inverses (for nonzero) elements. That is to say, given any x there is another element, namely $-x$, such that $x + (-x) = 0$, and given any nonzero x there is another element, namely $\frac{1}{x}$ such that $x \frac{1}{x} = 1$. Thus, a binary operation that has an identity element is said to have *inverses* if for each element there is an inverse element such that when combined by the operation they yield the identity element for the operation. Addition has inverses, and multiplication has inverses of nonzero elements.

Finally, there is a particular relation between the operations of addition and multiplication, that of distributivity:

$$x(y + z) = xy + xz \quad \text{and} \quad (y + z)x = yx + zx.$$

Multiplication distributes over addition, that is, when multiplying a sum by x we can distribute the x over the terms of the sum.

Exercise 1. On properties of operations.

(a). Is the binary operation $x * y = \frac{xy}{x + y}$ for positive x and y a commutative operation? That is, is it true that $x * y = y * x$ for all positive x and y ? Is it associative? Explain your answer.

(b). Is it true that $(w - x) - (y - z) = (w - y) - (x - z)$ is an identity for real numbers? Can you say why or why not? (The word "identity" is used for an equation which holds whenever both sides of the equation are defined and are equal.)

(c). Although multiplication in \mathbf{R} distributes over addition, addition doesn't distribute over multiplication. Give an example where it doesn't.

Algebraic structures. We'll define fields, rings, and groups as three kinds of algebraic structures. An algebraic structure will have an underlying set, binary operations, unary operations, and constants, that have some of the properties mentioned above like commutativity, associativity, identity elements, inverse elements, and distributivity. Different kinds of structures will have different operations and properties.

The algebraic structures are abstractions of familiar ones like those on the real numbers \mathbf{R} , but for each kind of structure there will be more than one example, as we'll see.

1.2.2 Fields

Informally, a field is a set equipped with four operations—addition, subtraction, multiplication, and division that have the usual properties. (They don't have to have the other operations that \mathbf{R} has, like powers, roots, logs, and the myriad other functions like $\sin x$.)

Definition 1.1 (Field). A *field* is a set equipped with two binary operations, one called *addition* and the other called *multiplication*, denoted in the usual manner, which are both commutative and associative, both have identity elements (the additive identity denoted 0 and the multiplicative identity denoted 1), addition has inverse elements (the inverse of x denoted $-x$), multiplication has inverses of nonzero elements (the inverse of x denoted $\frac{1}{x}$ or x^{-1}), multiplication distributes over addition, and $0 \neq 1$.

This definition will be spelled out in detail in chapter 2.

Of course, one example of a field is the field of real numbers \mathbf{R} . What are some others?

Example 1.2 (The field of rational numbers, \mathbf{Q}). Another example is the field of rational numbers. A rational number is the quotient of two integers a/b where the denominator is not 0. The set of all rational numbers is denoted \mathbf{Q} . We're familiar with the fact that the sum, difference, product, and quotient (when the denominator is not zero) of rational numbers is another rational number, so \mathbf{Q} has all the operations it needs to be a field, and since it's part of the field of the real numbers \mathbf{R} , its operations have the the properties necessary to be a field. We say that \mathbf{Q} is a *subfield* of \mathbf{R} and that \mathbf{R} is an *extension* of \mathbf{Q} . But \mathbf{Q} is not all of \mathbf{R} since there are irrational numbers like $\sqrt{2}$.

Example 1.3 (The field of complex numbers, \mathbf{C}). Yet another example is the field of complex numbers \mathbf{C} . A complex number is a number of the form $a + bi$ where a and b are real numbers and $i^2 = -1$. The field of real numbers \mathbf{R} is a subfield of \mathbf{C} . We'll review complex numbers before we use them. See *Dave's Short Course on Complex Numbers* at <http://www.clarku.edu/~djoyce/complex>

In chapter 2, we'll study fields in detail, and we'll look at many other fields. Some will only have a finite number of elements. (They won't be subfields of \mathbf{Q} .) Some will have \mathbf{Q} as a subfield but be subfields themselves of \mathbf{R} or \mathbf{C} . Some will be even larger.

Exercise 2. On fields. None of the following are fields. In each case, the operations of addition and multiplication are the usual ones.

- The integers \mathbf{Z} do not form a field. Why not?
- The positive real numbers $\{x \in \mathbf{R} \mid x > 0\}$ do not form a field. Why not?
- The set of real numbers between -10 and 10 , that is,

$$(-10, 10) = \{x \in \mathbf{R} \mid -10 < x < 10\}$$

is not a field. Why not?

1.2.3 Rings

Rings will have the three operations of addition, subtraction, and multiplication, but don't necessarily have division. Most of our rings will have commutative multiplication, but some won't, so we won't require that multiplication be commutative in our definition. All the rings we'll look at have a multiplicative identity, 1, so we'll include that in the definition.

Definition 1.4 (Ring). A *ring* is a set equipped with two binary operations, one called *addition* and the other called *multiplication*, denoted in the usual manner, which are both associative, addition is commutative, both have identity elements (the additive identity denoted 0 and the multiplicative identity denoted 1), addition has inverse elements (the inverse of x denoted $-x$), and multiplication distributes over addition. If multiplication is also commutative, then the ring is called a *commutative ring*.

Of course, all fields are automatically rings, in fact commutative rings, but what are some other rings?

Example 1.5 (The ring of integers, \mathbf{Z}). The ring of integers \mathbf{Z} includes all integers (whole numbers)—positive, negative, or 0. Addition, subtraction, and multiplication satisfy the requirements for a ring, indeed, a commutative ring. But there are no multiplicative inverses for any elements except 1 and -1 . For instance, $1/2$ is not an integer. We'll find that although the ring of integers looks like it has less structure than a field, this very lack of structure allows us to discover more about integers. We'll be able to talk about prime numbers, for example.

Example 1.6 (Polynomial rings). A whole family of examples are the rings of polynomials. Let R be any commutative ring (perhaps a field), and let $R[x]$ include all polynomials with coefficients in R . We know how to add, subtract, and multiply polynomials, and these operations have the properties required to make $R[x]$ a commutative ring. We have, for instance, the ring of polynomials with real coefficients $\mathbf{R}[x]$, the ring with integral coefficients $\mathbf{Z}[x]$, etc.

Example 1.7 (Matrix rings). How about an example ring that's not commutative? The ring of $n \times n$ matrices with entries in a commutative ring R gives such an example, this ring being denoted $M_n(R)$. This ring, $M_n(R)$, won't be commutative when $n \geq 2$. An example of a matrix ring is the ring of 2×2 matrices with real entries, $M_2(\mathbf{R})$. Addition and subtraction are computed coordinatewise. The additive identity, 0, of this matrix ring is the matrix with all 0 entries, $0 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$. Matrix multiplication is not coordinatewise, but it is associative, and multiplication does distribute over addition. The multiplicative identity for this matrix ring is what's usually called the identity matrix, denoted I . It has 1's down the main diagonal and 0's elsewhere, $1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$.

Sylvester (1814–1897), in 1850, called rectangular arrangements of numbers *matrices*, and Cayley wrote much about them in his papers of 1855–1858.

Example 1.8 (Integers modulo n). An important family of rings is the ring of integers modulo n . We'll study this in more detail later in section 2.2, but here's an incomplete overview. Fix a positive integer n . Think of two integers a and b as being the same modulo n if n divides $b - a$. In that case, we'll say that a and b are *congruent modulo n* , and we'll

use the notation Gauss (1777–1855) developed, $a \equiv b \pmod{n}$, to denote that congruence. Congruence is commonly used in the study of number theory. This meaning of the Latin word “modulo” was introduced into mathematics by Gauss in 1801.

Note that there are only n distinct integers modulo n , namely 0 through $n - 1$, since those are the only remainders you can get when you divide an integer by n . These remainders are also called “residues”. We can represent integers modulo n by these remainders from 0 through $n - 1$. Thus, we’ll say, for instance, that 5 plus 3 equals 1 modulo 7, by which we mean $5 + 3 \equiv 1 \pmod{7}$. Thus, we can turn congruence modulo n , which is an equivalence relation on \mathbf{Z} into equality on an n -element set. That n -element set is denoted $\mathbf{Z}/n\mathbf{Z}$, read \mathbf{Z} modulo $n\mathbf{Z}$, or more simply as \mathbf{Z}_n , read $\mathbf{Z} \bmod n$. So, we can take the elements of \mathbf{Z}_n to be the integers from 0 through $n - 1$, where we understand that addition, subtraction, and multiplication are done modulo n . And it turns out that this is a ring, as we’ll see when we study \mathbf{Z}_n in detail.

Incidentally, when n is a prime number p , then \mathbf{Z}_p is not just a ring, but a field, as will be discussed in section 2.2.

Exercise 3. On rings. None of the following are rings. In each case, the operations of addition and multiplication are the usual ones.

- (a). The set of nonzero integers, $\{x \in \mathbf{Z} \mid x \neq 0\}$ is not a ring. Why not?
- (b). The set of even integers $\{2x \mid x \in \mathbf{Z}\}$ is not a ring. Why not?
- (c). The set of odd degree polynomials with real coefficients

$$\{f(x) \in \mathbf{R}[x] \mid \text{the degree of } f(x) \text{ is odd}\}$$

is not a ring. Why not? (How about the set of even degree polynomials?)

Exercise 4. On noncommutative rings. Are the following rings? (The operations are the usual matrix operations.) Explain in a sentence or two, but a proof is not necessary.

- (a). The set of all matrices with real coefficients (all sizes).
- (b). The set of all 2×2 matrices with real entries of the form

$$\begin{bmatrix} a & b \\ 0 & d \end{bmatrix}.$$

- (c). The set of all 2×2 matrices with real entries of the form

$$\begin{bmatrix} a & b \\ -b & a \end{bmatrix}.$$

In chapter 3 we’ll analyze rings in more detail.

1.2.4 Groups

Unlike fields and rings which have two primary binary operations, groups only have one binary operation.

Definition 1.9 (Group). A *group* is a set equipped with a binary operation that is associative, has an identity element, and has inverse elements. If, furthermore, multiplication is

commutative, then the group is called a *commutative group* or an *Abelian group*. Abelian groups can be denoted either additively or multiplicatively, but nonabelian groups are usually denoted multiplicatively. We'll use the term *order of the group* to indicate how many elements a group G has and denote this order by $|G|$.

Example 1.10 (The underlying additive group of a ring). Of course, if you have a field or ring, and just consider addition (and forget about multiplication) you've got an Abelian group. Sometimes this is called the *underlying additive group* of the field or ring. We'll use the same notation for the underlying additive group as we do for the ring. Thus, \mathbf{Z} could mean either the ring of integers or the Abelian group of integers under addition, depending on the context.

Example 1.11 (Finite cyclic groups). The underlying group of the ring \mathbf{Z}_n is called a cyclic group. Its elements are, of course, $0, 1, 2, \dots, n-1$ where n is congruent to 0. Cyclic groups are also written multiplicatively, and then the elements are $1, a, a^2, \dots, a^{n-1}$ where $a^n = 1$. A common notation for this cyclic group is C_n .

Definition 1.12 (Units in a ring). In order to use the multiplication for a group operation, we'll have to only include the units, also called invertible elements. A *unit* or *invertible element* of a ring R is an element $x \in R$ such that there exists another element $y \in R$ so that $xy = yx = 1$. The subset of units is denoted

$$R^* = \{x \in R \mid \exists y \in R, xy = 1\}.$$

You can easily show that the units form a group under multiplication, called the *multiplicative group of units* of R . When R is a field, then R^* is all of R except 0, but for rings there will be other elements than 0 that aren't invertible. The group R^* will be Abelian when the ring R is commutative, but usually it will be nonabelian when R is not commutative.

Examples 1.13. The units in the ring \mathbf{Z} are just 1 and -1 . The group of units \mathbf{Z}^* is a cyclic group of order 2.

We'll see later that the group of units Z_p^* when p is prime is a cyclic group of order $p-1$. It is usually the case that Z_n^* when n is composite is not a cyclic group.

Example 1.14 (A general linear group, $GL_2(\mathbf{R})$). As a particular example of a multiplicative group of units, take the invertible elements of the matrix ring $M_2(\mathbf{R})$. The invertible 2×2 matrices are those matrices

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

whose determinants $ad - bc$ are nonzero. The group of invertible $n \times n$ matrices, $M_n(R)^*$, is the *general linear group* with coefficients in the ring R , denoted $GL_n(R)$. Note that $GL_n(R)$ is a nonabelian group for $n \geq 2$. The real general linear group $GL_2(\mathbf{R})$ can be interpreted as the group of invertible linear transformations of the plane \mathbf{R}^2 that leave the origin fixed.

We'll study $GL_2(\mathbf{R})$ and $GL_n(\mathbf{R})$ in more detail in section 4.7.2.

Exercise 5. Find two matrices in $GL_2(\mathbf{Z})$ that don't commute thereby proving $GL_2(\mathbf{Z})$ is a nonabelian group.

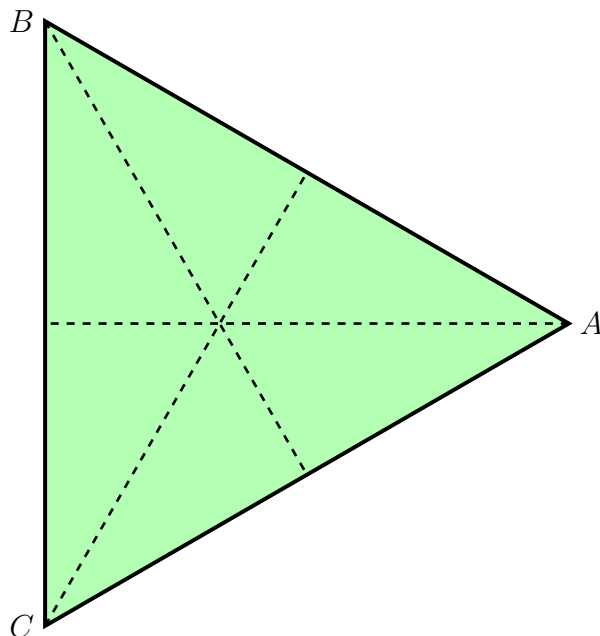


Figure 1.1: Equilateral triangle with lines of symmetry

There are many examples of finite nonabelian groups found in geometry. We'll look at the group of symmetries of an equilateral triangle.

Example 1.15 (The dihedral group D_3). Consider an equilateral triangle. Place a coordinate system on the plane of the triangle so that its center is at $(0, 0)$, one vertex, A , at $(1, 0)$, and the other two, B and C , at $(-\frac{1}{2}, \pm\frac{1}{2}\sqrt{3})$. This triangle has six symmetries. A symmetry is a transformation of the plane that preserves distance (that is, an *isometry*) that maps the triangle back to itself. Three of these symmetries are rotations by 0° , 120° , and 240° .

$$1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \rho = \begin{bmatrix} -\frac{1}{2} & -\frac{1}{2}\sqrt{3} \\ \frac{1}{2}\sqrt{3} & -\frac{1}{2} \end{bmatrix} \quad \rho^2 = \begin{bmatrix} -\frac{1}{2} & \frac{1}{2}\sqrt{3} \\ -\frac{1}{2}\sqrt{3} & -\frac{1}{2} \end{bmatrix}$$

The identity transformation, 1 , fixes A , B , and C ; the rotation ρ by 120° maps A to B , B to C , and C to A ; and the rotation ρ^2 by 240° maps A to C , B to A , and C to B . There are also three reflections.

$$\varphi = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \quad \rho\varphi = \begin{bmatrix} -\frac{1}{2} & \frac{1}{2}\sqrt{3} \\ -\frac{1}{2}\sqrt{3} & \frac{1}{2} \end{bmatrix} \quad \rho^2\varphi = \begin{bmatrix} -\frac{1}{2} & \frac{1}{2}\sqrt{3} \\ \frac{1}{2}\sqrt{3} & \frac{1}{2} \end{bmatrix}$$

The reflection φ fixes A , and interchanges B and C ; the reflection $\rho\varphi$ fixes C and interchanges A and B ; and the reflection $\rho^2\varphi$ fixes B and interchanges A and C . This is a particular nonabelian group that has 6 elements. It is a subgroup of $GL_2(\mathbf{R})$ mentioned above.

Example 1.16 (A group of functions). Many applications of group theory are to groups of invertible functions. Such a group includes invertible functions on some set such that the composition of any two of the functions is another one.

Let $f(x) = 1/x$ and $g(x) = 1 - x$. Both of those are invertible considered as rational functions, and, in fact, each is its own inverse: $(f \circ f)(x) = f(1/x) = 1$, and $(g \circ g)(x) =$

$g(1-x) = 1 - (1-x) = x$. Let's see what other functions we can derive from f and g by composing them.

First, consider $(f \circ g)(x) = f(g(x)) = f(1-x) = \frac{1}{1-x}$; call that composition h so that $h(x) = \frac{1}{1-x}$. Next, consider $(g \circ f)(x) = g(f(x)) = g\left(\frac{1}{x}\right) = 1 - \frac{1}{x} = \frac{x-1}{x}$; call that composition k so that $k(x) = \frac{x-1}{x}$.

We can get more functions if we continue to compose these. Note that $(f \circ k)(x) = f\left(\frac{x-1}{x}\right) = \frac{x}{x-1}$; call that ℓ so that $\ell(x) = \frac{x}{x-1}$. Also, $(g \circ h)(x) = g\left(\frac{1}{1-x}\right) = 1 - \frac{1}{1-x} = \frac{x}{x-1}$. That function has already been called ℓ , so $g \circ h = \ell$.

A couple more computations show that $h \circ h = k$ and $k \circ k = h$.

Since f and g are each their own inverses, $f \circ f = i$ and $g \circ g = i$, where i is the identity function, $i(x) = x$. Also $h \circ k = k \circ h = i$, and $\ell \circ \ell = i$. Also, i composed with any function (on either side) is equal to that same function.

It turns out that these six functions are closed under composition. Table 1.1 gives all of their compositions.

	i	f	g	h	k	ℓ
i	i	f	g	h	k	ℓ
f	f	i	h	g	ℓ	k
g	g	k	i	ℓ	f	h
h	h	ℓ	f	k	i	g
k	k	g	ℓ	i	h	f
ℓ	ℓ	h	k	f	g	i

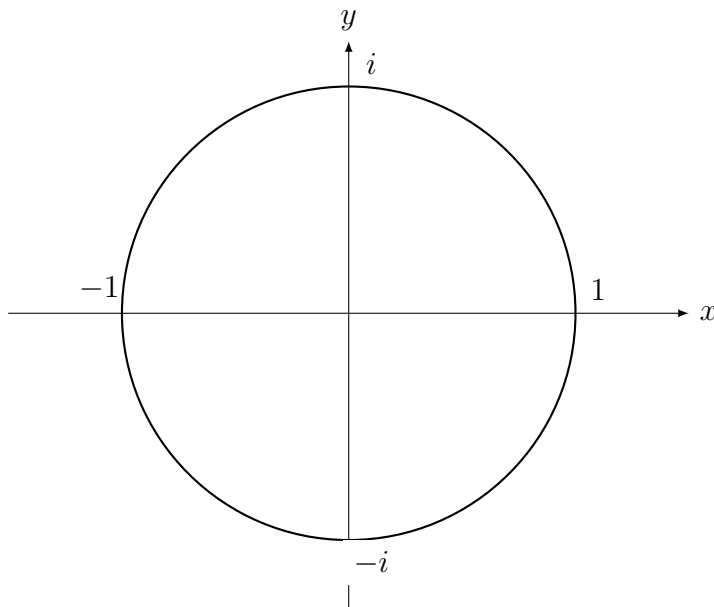
Table 1.1: Composition table for six particular rational functions.

Note that in each row and each column of the table, each one of the functions appears exactly once. That makes the entries of the table a Latin square. A *Latin square* is a square $n \times n$ array filled with n different symbols, each occurring exactly once in each row and exactly once in each column.

Example 1.17 (Euler's circle group). The unit circle, $S^1 = \{x + yi \in \mathbf{C} \mid x^2 + y^2 = 1\}$, is a group under multiplication. This is sometimes called Euler's circle group since Euler (1707–1783) introduced the unit circle in the complex plane for studying angles and trigonometric functions.

The product of two complex numbers on this unit circle is another number on the unit circle. You can directly verify that or you can show it by trigonometry. If $x + yi$ is on the unit circle, then we can identify x with $\cos \theta$ and y with $\sin \theta$ where θ is, as usual, the angle between the positive x -axis and the ray from 0 to $x + yi$. Then the product of two complex numbers on the unit circle corresponds to adding their angles together. The addition formulas for cosines and sines give this correspondence.

Exercise 6. Compute the product of $\cos \theta + i \sin \theta$ times $\cos \varphi + i \sin \varphi$. If $x + iy = (\cos \theta + i \sin \theta)(\cos \varphi + i \sin \varphi)$, then what is x , the real part of the product, in terms of θ and φ ? What is y , the imaginary part?

Figure 1.2: Unit circle S^1

Comment 1.18. Although the sphere

$$S^2 = \{(x, y, z) \in \mathbf{R}^3 \mid x^2 + y^2 + z^2 = 1\}$$

has no group structure, the 3-sphere in 4-space does. The 3-sphere is

$$S^3 = \{(x, y, z, w) \in \mathbf{R}^4 \mid x^2 + y^2 + z^2 + w^2 = 1\}.$$

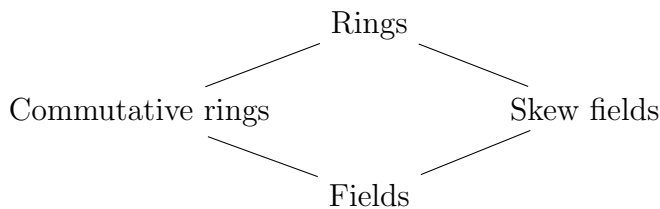
We don't have time or space to discuss that group structure here. (The 2-sphere S^2 , in fact, spheres in all dimensions, does have quandle structures, whatever a quandle might be. See section 4.5.2.)

In chapter 4 we'll study groups in detail.

1.2.5 Other algebraic structures besides fields, rings, and groups

There are an unlimited number of other algebraic structures. Some are similar to those listed above.

For instance, there are division rings (also called skew fields) that have all the properties of fields except multiplication doesn't have to be commutative. The primary example is the quaternions \mathbf{H} . We'll discuss quaternions later in section 2.5.2.



There are a number of structures that are just commutative rings that have nice properties, and we'll look at some of them including integral domains, unique factorization domains, principal ideal domains, and Euclidean domains.

Sometimes rings that don't have a multiplicative identity are studied, but for us, we'll always have 1.

You've already studied vector spaces over the real numbers. Most of the things that you've studied about vector spaces over \mathbf{R} also hold for vector spaces over other fields.

The analogous structure for vector spaces when a field is replaced by a ring is called a *module* over the ring. We won't study modules over a ring, but when we look at ideals in a ring, they are, in fact, examples of modules over the ring. Also, Abelian groups are modules over the ring \mathbf{Z} .

We'll discuss another algebraic structure, quandles, in section 4.5.2 when we discuss groups.

1.3 Isomorphisms, homomorphisms, etc.

Frequently, we look at two algebraic structures A and B of the same kind, for instance, two groups or two rings or two fields, and we'll want to compare them. For instance, we might think they're really the same thing, but they have different names for their elements. That leads to the concept of isomorphism $f : A \cong B$, and we'll talk about that first. Other times we'll know they're not the same thing, but there is a relation between them, and that will lead to the next concept, homomorphism, $f : A \rightarrow B$. We'll then look at some special homomorphisms such as monomorphisms. When we have a homomorphism $f : A \rightarrow A$, we'll call it an endomorphism, and when an isomorphism $f : A \cong A$, we'll call it an automorphism. We'll take each of these variants in turn.

The concepts of injection (one-to-one function), surjection (onto function), and bijection are described section A.2.2 in the appendix on functions.

We'll use the following theorem about finite sets when we consider homomorphisms between finite algebraic structures.

Theorem 1.19. Suppose that $f : A \rightarrow B$ is a function between two finite sets of the same cardinality. Then the following three conditions are equivalent: (1) f is a bijection, (2) f is an injection, and (3) f is a surjection.

Exercise 7. Prove that if $f : A \rightarrow B$ is a function between two finite sets of the same cardinality, then f is injective if and only if f is surjective.

1.3.1 Isomorphisms

We'll say two algebraic structures A and B are isomorphic if they have exactly the same structure, but their elements may be different. For instance, let A be the ring $\mathbf{R}[x]$ of polynomials in the variable x with real coefficients while B is the ring $\mathbf{R}[y]$ of polynomials in y . They're both just polynomials in one variable, it's just that the choice of variable is different in the two rings. We need to make this concept more precise.

Definition 1.20 (Ring isomorphism). Two rings A and B are *isomorphic* if there is a bijection $f : A \rightarrow B$ which preserves addition and multiplication, that is, for all x and y in A ,

$$f(x + y) = f(x) + f(y), \text{ and } f(xy) = f(x)f(y).$$

The correspondence f is called a *ring isomorphism*.

After we introduce homomorphisms, we'll have another way to describe isomorphisms.

You can prove various properties of ring isomorphism from this definition.

Exercise 8. Since the structure of rings is defined in terms of addition and multiplication, if f is a ring isomorphism, it will preserve structure defined in terms of them. Verify that f preserves 0, 1, negation, and subtraction.

Exercise 9. Prove that if f is a ring isomorphism, then so is its inverse function $f^{-1} : B \rightarrow A$.

Exercise 10. Prove that if $f : A \rightarrow B$ and $g : B \rightarrow C$ are both ring isomorphisms, then so is their composition $(g \circ f) : A \rightarrow C$.

Since a field is a special kind of ring, and its structure is defined in terms of addition and multiplication, we don't need a special definition for a field isomorphism. A field isomorphism is just a ring isomorphism between fields.

Exercise 11. Prove that if a ring is isomorphic to a field, then that ring is a field.

We do need a different definition for a group isomorphism since a group is defined in terms of just one binary operation instead of two.

Definition 1.21 (Group isomorphism). Two groups A and B are isomorphic if there is a bijection $f : A \rightarrow B$ which preserves the binary operation. If both are written additively, that means for all x and y in A , $f(x + y) = f(x) + f(y)$; if multiplicative notation is used in both, then $f(xy) = f(x)f(y)$; if additive in A but multiplicative in B , then $f(x + y) = f(x)f(y)$; and if multiplicative in A and additive in B , then $f(xy) = f(x) + f(y)$. The correspondence f is called a *group isomorphism*.

Usually A and B will use the same notation, both additive or both multiplicative, but not always.

Exercise 12. Suppose that both A and B are written multiplicatively and that $f : A \rightarrow B$ is a group isomorphism. Prove that $f(1) = 1$ and $f(x^{-1}) = f(x)^{-1}$ for all $x \in A$.

Example 1.22. Let $A = \mathbf{Z}$ be the group of integers under addition. Let B be the integral powers of 2, so $B = \{\dots, \frac{1}{4}, \frac{1}{2}, 1, 2, 4, \dots\}$ with multiplication as the operation in B . Prove that an isomorphism $f : A \rightarrow B$ is defined by $f(n) = 2^n$. You'll need to show that $f(m+n) = f(m)f(n)$.

There's actually another isomorphism $g : A \rightarrow B$, too, defined by $g(n) = 2^{-n}$.

1.3.2 Homomorphisms

Whereas isomorphisms are bijections that preserve the algebraic structure, homomorphisms are simply functions that preserve the algebraic structure. Since the word homomorphism is so long, alternate words are often used like morphism and map, especially in spoken mathematics.

Definition 1.23 (Ring homomorphism). A *ring homomorphism* $f : A \rightarrow B$ between rings is a function that preserves addition, multiplication, and 1.

A *group homomorphism* $f : A \rightarrow B$ between groups preserves the binary operation (addition or multiplication depending on the notation used for the group).

Comment 1.24. It's a peculiarity of rings that preserving addition and multiplication doesn't imply that 1 is also preserved, so that condition has to be required as well. We'll see plenty of examples of homomorphisms in the course, and there are more examples in the next section on monomorphisms. Of course, isomorphisms are special cases of homomorphisms.

Example 1.25 (A ring homomorphism). Let $\mathbf{Z}[x]$ be the ring of polynomials with integral coefficients. Evaluating a polynomial $f(x)$ at a particular number, like 3, to give $f(3)$, is a ring homomorphism $\varphi : \mathbf{Z}[x] \rightarrow \mathbf{Z}$. It preserves addition since $\varphi(f(x) + g(x)) = f(3) + g(3) = \varphi(f(x)) + \varphi(g(x))$, and you can check that it preserves multiplication and 1.

Example 1.26 (A group homomorphism). Let A be the integers under addition, and let $B = \{1, -1\}$ with multiplication as the binary operation. Then $f : A \rightarrow B$ defined by $f(n) = (-1)^n$ is a group homomorphism.

You can prove several properties of homomorphisms from the definition, but for the time being I'll just mention two because they'll lead to the concept of *category* which will be introduced in section 3.5.

1. The composition of two homomorphisms (of the same kind) is another homomorphism.
2. The identity function $1_A : A \rightarrow A$, which maps every element to itself, is a homomorphism, indeed, it's an isomorphism.

When we have a homomorphism $f : A \rightarrow B$, we'll call A the *domain* of f and we'll call B the *codomain* of f . (Sometimes the word "range" is used for codomain, but some people prefer to use "range" to mean image, which is a different thing. To avoid ambiguity, we'll use "codomain".)

A more natural way to characterize isomorphism is in terms of homomorphisms. Two rings A and B are isomorphic if and only if there are two ring homomorphisms $f : A \rightarrow B$ and $g : B \rightarrow A$ such that $g \circ f$ is the identity on A and $f \circ g$ is the identity on B .

1.3.3 Monomorphisms and epimorphisms

Two common kinds of homomorphisms are monomorphisms and epimorphisms, often called monos and epis for short. When a homomorphism $f : A \rightarrow B$ is an injective function, it's called a *monomorphism*; and when it is a surjective function, it's an *epimorphism* (but, in the category of rings, we'll see there are more epimorphisms than just the surjective ring homomorphisms). You might wonder why we need these words when we've got more than enough words already to describe injective (one-to-one) and surjective (onto) as well as others not mentioned here. The main reason is that they're special kinds of injections or surjections—they preserve the algebraic structure. Another is that, although for group homomorphisms monos and epis have these particular correspondences to injective and surjective, there are other categories in which they don't.

Note that every isomorphism is simultaneously a monomorphism and an epimorphism. The converse holds for groups, but, surprisingly, not for rings.

Example 1.27 (Inclusion). Inclusions are monomorphisms. When one ring (or group) A is a subring (or subgroup) of another B , then the inclusion function $\iota : A \rightarrow B$, which maps an element to itself, is a monomorphism. That's an important example of a monomorphism, but there are others.

Example 1.28. For example, let A and B both be the additive group of integers \mathbf{Z} , and let $f(n) = 2n$. This f is a monomorphism, but it's not an inclusion (which in this case would be the identity map since A and B are the same).

Comment 1.29. Note that if $f : A \rightarrow B$ is a ring homomorphism where A is a field and $0 \neq 1$ in B , then f is always an injection, and so it's a monomorphism. You can prove this statement in two stages. First, show that if $f(x) = 0$ then $x = 0$. Second, show that if $f(x) = f(y)$, then $x = y$.

Thus, every field homomorphism is a monomorphism.

Example 1.30 (A group epimorphism). We'll see plenty of epimorphisms when we talk more about the integers modulo n , but for the time being, consider example 1.26 of a group epimorphism. The group A is the additive group of integers \mathbf{Z} , and the group B is the two element group $\{1, -1\}$ under multiplication. Then $f : A \rightarrow B$ defined by $f(n) = (-1)^n$ is a group epimorphism. Even numbers are sent to 1 and odd numbers to -1 .

1.3.4 Endomorphisms and automorphisms

An endomorphism is just a homomorphism $f : A \rightarrow A$ where the domain and codomain are the same, and an automorphism is just an isomorphism $f : A \rightarrow A$. These are important because we always have the identity automorphism $1_A : A \rightarrow A$ to compare f to, so we have more information when the domain and codomain are the same.

Example 1.31 (A field automorphism). Let \mathbf{C} be the complex field. Let $\phi : \mathbf{C} \rightarrow \mathbf{C}$ be *complex conjugation*, usually denoted by putting a bar above the complex number

$$\varphi(x + yi) = \overline{x + yi} = x - yi.$$

This is clearly a bijection since it is its own inverse, $\overline{\overline{x + yi}} = x + yi$. Also, it preserves addition, multiplication, and 1, so it's a ring isomorphism.

$$\begin{aligned} \overline{(x_1 + y_1i) + (x_2 + y_2i)} &= \overline{x_1 + y_1i} + \overline{x_2 + y_2i} \\ \overline{(x_1 + y_1i)(x_2 + y_2i)} &= \overline{x_1 + y_1i} \overline{x_2 + y_2i} \\ \overline{1} &= 1 \end{aligned}$$

In fact, it's a field automorphism of \mathbf{C} .

The existence of this automorphism says that we can't distinguish between i and $-i$ in the sense that any true statement about the complex numbers remains true when all occurrences of i are replaced by $-i$.

Example 1.32 (Group endomorphisms and automorphisms). There are many group endomorphisms $f : \mathbf{Z} \rightarrow \mathbf{Z}$ from the additive group of integers to itself. Fix any integer n and let $f(x) = nx$. This is a group homomorphism since $f(x+y) = n(x+y) = nx + ny = f(x) + f(y)$.

For $n \neq 0$ it is also a monomorphism. For $n = -1$ this is negation, and it's a bijection, so it's a group automorphism. That says if we only consider addition, we can't distinguish between positive and negative numbers.

But negation is not a ring automorphism on the ring of integers because $-(xy)$ does not equal $(-x)(-y)$. Thus, with the use of multiplication, we can distinguish between positive and negative numbers.

1.4 A little number theory

In science nothing capable of proof ought to be accepted without proof. Though this demand seems so reasonable, yet I cannot regard it as having been met even in the most recent methods of laying the foundations for the simplest science; viz., that part of logic which deals with the theory of numbers.

Dedekind, 1888

This course is not meant to be a course in number theory, but we will need a little bit of it. We'll quickly review mathematical induction on the natural numbers \mathbf{N} , divisibility, prime numbers, greatest common divisors, and the Euclidean algorithm.

1.4.1 Mathematical induction on the natural numbers \mathbf{N}

Richard Dedekind (1831–1916) published in 1888 a paper entitled *Was sind und was sollen die Zahlen?* variously translated as *What are numbers and what should they be?* or *The Nature of Meaning of Numbers*. In that work he developed basic set theory and characterized the natural numbers as a simply infinite set.

Definition 1.33. (Dedekind) A set \mathbf{N} is said to be *simply infinite* when there exists a one-to-one function $\mathbf{N} \xrightarrow{\prime} \mathbf{N}$ called the *successor function*, such that there is an element, called the *initial element* and denoted 1, that is not the successor of any element, and if a subset S of \mathbf{N} contains 1 and is closed under the successor function, then $S = \mathbf{N}$.

Such a simply infinite set \mathbf{N} may be called the *natural numbers*. It is characterized by an element 1 and a transformation $\mathbf{N} \xrightarrow{\prime} \mathbf{N}$ satisfying the following conditions:

1. Injectivity: $\forall n, m, n \neq m$ implies $n' \neq m'$.
2. Initial element: $\forall n, 1 \neq n'$.
3. Induction: If $S \subseteq \mathbf{N}$, $1 \in S$, and $(\forall n, n \in S$ implies $n' \in S)$, then $S = \mathbf{N}$.

The Dedekind axioms, also called the Peano axioms, are this last characterization involving 1, the successor function, and the three conditions. Among other things, Peano (1858–1932) developed much of the notation in common use in set theory.

The last axiom is called mathematical induction. If you want to show a subset S of \mathbf{N} is all of \mathbf{N} , first show that $1 \in S$. Then show for each natural number n that $n \in S$ implies $n + 1$ in S . Finally conclude that $S = \mathbf{N}$.

A principle that is logically equivalent to mathematical induction is the well-ordering principle, also called the minimization principle. It says that each nonempty subset of \mathbf{N} has a least element. To use it to prove a subset S of \mathbf{N} is all of \mathbf{N} , assume that it isn't, take the least element n in $\mathbf{N} - S$, and derive a contradiction, usually by showing there's a smaller element than n not in S .

Another principle logically equivalent to mathematical induction is Euclid's principle of infinite descent which says that there is no infinite decreasing sequence of positive integers. This principle was also used by Fermat (1607–1665).

1.4.2 Divisibility

We'll restrict our discussion now to \mathbf{N} , the natural numbers, that is, the set of positive integers.

Recall that an integer m *divides* an integer n , written $m|n$, if there exists an integer k such that $mk = n$. A few basic properties of divisibility follow directly from this definition. Euclid (fl. ca. 300 B.C.E.) uses some of these in Book VII of his *Elements*. You can find Joyce's translation of Euclid's *Elements* on the web at <http://aleph0.clarku.edu/~djoyce/java/elements/elements.html>

1. 1 divides every number. $1|n$.
2. Each number divides itself. $n|n$.
3. If one number m divides another number n , then m divides any multiple of n , that is, $m|n$ implies $m|kn$.
4. Divisibility is a transitive relation, that is, $m|n$ and $n|k$ imply $m|k$.
5. If one number divides two other numbers, then it divides both their sum and difference. $m|n$ and $m|k$ imply $m|(n + k)$ and $m|(n - k)$.
6. Cancellation law. One number divides another if and only if any multiple of that one number divides the same nonzero multiple of the other number. $m|n \iff kn|kn$. ($k \neq 0$)

Example 1.34. The divisors of a number can be displayed graphically in what is called a Hasse diagram of the lattice of divisors. As an example, consider the number 432. Its prime factorization is 2^43^3 , so its divisors are of the form 2^m3^n where $0 \leq m \leq 4$ and $0 \leq n \leq 3$. There are $5 \cdot 4 = 20$ of these divisors. They are

1	2	3	4	6	8	9	12	16	18
24	27	36	48	54	72	108	144	216	432

We can display these numbers and emphasize which ones divide which other ones if we put the large numbers at the top of the diagram, and connect the smaller divisors to the larger ones with lines. That results in the Hasse diagram in figure 1.3.

Since divisibility is transitive, we don't have to include all possible connections. So long as there is a path of connections from a lower number to an upper one, then we can conclude the lower divides the upper. The resulting diagram is called a *Hasse diagram* in honor of Hasse (1898–1979).

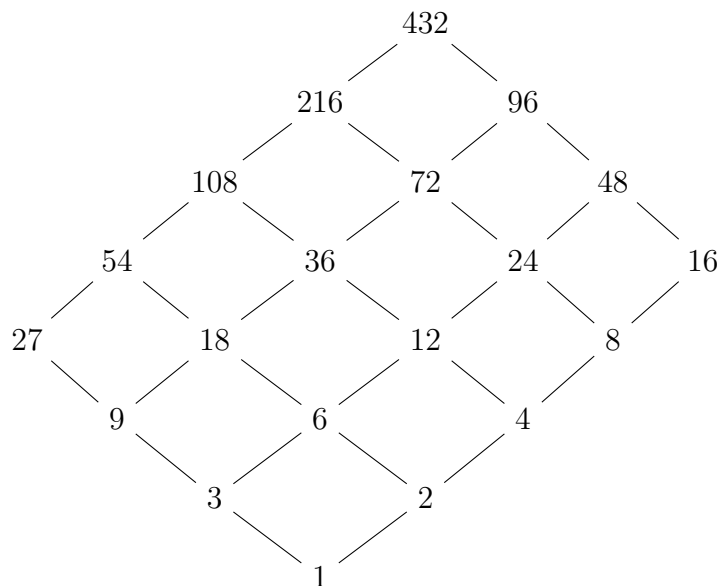


Figure 1.3: Divisors of 432

Exercise 13. Draw Hasse diagrams for the divisors of 30, 32, and 60.

The Hasse diagram for all positive integers under divisibility is, of course, infinite. Figure 1.4 shows the part of it up through 12.

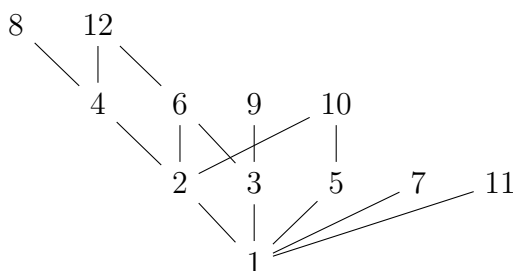


Figure 1.4: Divisibility up through 12

1.4.3 Prime numbers

Definition 1.35. A natural number greater than 1 is said to be a *prime number*, or more simply a *prime*, if its only divisors are 1 and itself, but if it has more divisors, it's called a *composite number*.

Two positive integers are said to be *relatively prime*, or *coprime* if the only positive integer that divides them both is 1.

Prime numbers were mentioned by the Pythagoreans Philolaus (470–385 B.C.E.) and Thymaridas (400–350 B.C.E.), and by Aristotle (384–322 B.C.E.) after them. The first recorded proofs about prime numbers occur in Euclid's *Elements*.

We know intuitively that there are infinitely many primes, and that every number is a product of primes. Now let's prove those statements. We'll start by proving something that will help us prove these two statements. If a theorem is not particularly interesting, but is useful in proving an interesting statement, then it's often called a lemma. This one is found in Euclid's *Elements*.

Lemma 1.36 (Euclid, VII.31). Every number greater than 1 has at least one prime divisor.

Proof. Let n be an integer greater than 1. We'll find a prime divisor of n . Let m be the smallest divisor of n greater than 1. (Note that we're using the minimization principle, also called the well-ordering principle, to conclude that such an m exists.) We'll show that m is prime thereby proving the lemma. We'll do that with a proof by contradiction, and that means that first we'll suppose that m is not prime, then derive a contradiction, and that will imply that m must be prime.

Suppose m is not prime, but composite. Then m is the product of two integers, j and k , each greater than 1. Now, $k|m$ and $m|n$, so $k|n$. But $k < m$. That gives us a divisor of n which is even smaller than m but still greater than 1. That contradicts the fact that m is the smallest divisor of n greater than 1. Thus, m is prime, and it's a divisor of n . Q.E.D.

Now we can prove one of the two statements.

Theorem 1.37. Every number greater than 1 is either a prime or the product of primes.

Proof. This will be another proof by contradiction that uses the well-ordering principle.

Suppose that the theorem is false. Then there is some composite number greater than 1 that is not the product of primes. Let n be the smallest such. By our lemma, this n has some prime divisor, call it p . Then $m = n/p$ is a number smaller than n but larger than 1, so, by the minimality of n , m is either prime or the product of primes. In the first case, when m is prime, then $n = pm$ is the product of two primes. In the second case when m is a product of primes, then $n = pm$ is also a product of primes. In any case, n is the product of primes, a contradiction. Thus, the theorem is true. Q.E.D.

This last theorem will form part of the so-called fundamental theorem of arithmetic that says every number greater than 1 can be uniquely factored as a product of primes. So far, we only have that every number is a product of primes, but we haven't seen the uniqueness. We'll prove that pretty soon.

Next, let's prove the other statement, that there are infinitely many primes. This is Euclid's proof.

Theorem 1.38 (Euclid IX.20). There are infinitely many primes.

Proof. Actually, Euclid proves something a little stronger. Given any finite list of primes, he finds a prime not on that list.

Suppose that p_1, p_2, \dots, p_k is a finite list of primes. Let n be the product of these primes,

$$n = p_1 p_2 \cdots p_k.$$

By our lemma $n + 1$ has a prime factor, call it p . This prime p cannot equal any p_i , for then p would divide both n and $n + 1$, and so would divide the difference 1. But a prime p can't divide 1 since $p > 1$. This p is a prime not on the list.

It follows that there are infinitely many primes. Q.E.D.

The number of relatively prime integers. An important combinatorial count for number theory and algebra is the number $\varphi(n)$ of positive integers less than a given integer n . For example, we'll show later in corollary 2.11 that the number of units in the ring \mathbf{Z}_n is $\varphi(n)$. We'll also use it in our discussion of cyclotomic polynomials in section 1.6.2.

It's easy enough to compute $\varphi(n)$ when n is small. For example, $\varphi(12) = 4$, since there are four positive integers less than 12 which are relatively prime to 12, namely, 1, 5, 7, and 11.

Definition 1.39 (Euler's totient function). For a given positive integer n , the number of positive integers less than n that are relatively prime to n is denoted $\varphi(n)$. The function φ is called Euler's totient function.

The first few values of the totient function are listed in this table.

n	1	1	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
$\varphi(n)$	1	1	2	2	4	2	6	4	6	4	10	4	12	6	8	8	16	6

One obvious property of this function is that if p is prime, then $\varphi(p) = p - 1$.

A property that's not so obvious is that if m and n are relatively prime, then $\varphi(mn) = \varphi(m)\varphi(n)$. That property is summarized by saying that φ is a *multiplicative* function. It follows from the Chinese remainder theorem discussed in section 3.2.1.

That reduces the computation of φ to computing it on powers p^k of prime numbers. That can be found directly. The only positive integers less than or equal to p^k that aren't relatively prime to p^k are the multiples of p , which are $p, 2p, \dots, p^k$, and there are p^{k-1} of them. Therefore, $\varphi(p^k) = p^k - p^{k-1} = p^k \left(1 - \frac{1}{p}\right)$.

Theorem 1.40 (Euler's product formula).

$$\varphi(n) = n \prod_{p|n} \left(1 - \frac{1}{p}\right).$$

Proof. Write $n = p_1^{k_1} \cdots p_r^{k_r}$ as a product of powers of distinct primes. Then by the multiplicativity of φ ,

$$\varphi(n) = \varphi(p_1^{k_1}) \cdots \varphi(p_r^{k_r}) = p_1^{k_1} \left(1 - \frac{1}{p_1}\right) \cdots p_r^{k_r} \left(1 - \frac{1}{p_r}\right) = \varphi(n) = n \prod_{p|n} \left(1 - \frac{1}{p}\right).$$

Q.E.D.

1.4.4 The Euclidean algorithm

The Euclidean algorithm is an algorithm to compute the greatest common divisor of two natural numbers m and n . Euclid described in Book VII of his *Elements*.

Euclid defined the *greatest common divisor* of two natural numbers m and n , often denoted $\text{GCD}(m, n)$ or more simply just (m, n) , as the largest number d which is at the same time a divisor of m and a divisor of n .

Among other things, greatest common divisors are used to reduce common fractions to lowest terms. For example, if you wanted to reduce the fraction $\frac{1417}{1853}$ to lowest terms, you would look for the greatest common divisor of the two numbers 1417 and 1853, which is 109. Then you could divide both the numerator and the denominator by that greatest common divisor to reduce $\frac{1417}{1853}$ to its lowest terms, namely, $\frac{13}{17}$.

There are two forms of the Euclidean algorithm. The first form, as Euclid stated it, repeatedly subtracts the smaller number from the larger replacing the larger by the difference, until the two numbers are reduced to the same number, and that's the greatest common divisor. (Note that the process has to stop by the well-ordering principle since at each step the larger number is reduced.)

The other form speeds up the process. Repeatedly divide the smaller number into the larger replacing the larger by the remainder. (This speeds up the process because if the smaller number is much smaller than the larger, you don't have to subtract it from the larger many times, just divide once and take the remainder which is the same as what you'd get if repeatedly subtracted it.)

Example 1.41. Let's find $\text{GCD}(6731, 5777)$. Since $6731 - 5777 = 954$, replace 6731 by 954. We've reduced the problem to finding $\text{GCD}(5777, 954)$.

Now repeatedly subtract 954 from 5777 until you get a number smaller than 954 and replace 5777 by that number. Alternatively, you could divide 954 into 5777 and replace 5777 by the remainder. You'll get the same thing, namely 53.

Next to find $\text{GCD}(954, 53)$. If you keep subtracting 53 from 954, eventually you'll get 0. Or if you're using division, when you divide 53 into 954, you'll get a remainder of 0. Either way, you can conclude 53 divides 954, so their GCD is 53 itself. Thus, $\text{GCD}(6731, 5777) = 53$

This Euclidean algorithm works to produce the GCD, and the argument only depended on two properties of divisibility mentioned above, namely that if one number divides two other numbers, then it divides both their sum and difference.

Sometimes the GCD of two numbers turns out to be 1, and in that case we say the two numbers are *relatively prime* or that they're *coprime*.

Theorem 1.42 (Euclidean algorithm). Let d be the result of applying the Euclidean algorithm to m and n . Then d is the greatest common divisor $\text{GCD}(m, n)$. Furthermore, the common divisors k of m and n are the divisors of $\text{GCD}(m, n)$.

Proof. One step of the Euclidean algorithm replaces the pair (m, n) by $(m - n, n)$. It was mentioned above in the properties of divisibility that if one number divides two other numbers, then it divides both their sum and difference. Therefore, a number k divides both m and n if and only if k divides $m - n$ and n . Since the pair (m, n) have the same set of divisors as the pair $(m - n, n)$, therefore $\text{GCD}(m, n) = \text{GCD}(m - n, n)$. Thus, at each step of the Euclidean algorithm the GCD remains invariant. Eventually, the two numbers are the same, but when that last step is reached, that number is the GCD. So, the end result of the Euclidean algorithm is $d = \text{GCD}(m, n)$.

The remarks above show that every divisor k of m and n also divides the result d of applying the Euclidean algorithm to m and n . Finally, if $k|d$, since $d|m$ and $d|n$, therefore $k|m$ and $k|n$. Q.E.D.

Extended Euclidean algorithm. There's still more that we can get out of the algorithm if we include the equations implicit in the computations. That will lead to the extended Euclidean algorithm.

Example 1.43. When we found $\text{GCD}(6731, 5777)$, if we kept track of the quotients as well as the remainders, then each step yields an equation.

$$\begin{aligned} 6731 - 1 \cdot 5777 &= 954 \\ 5777 - 6 \cdot 954 &= 53 \\ 954 - 18 \cdot 53 &= 0 \end{aligned}$$

Turning these equations around, we can find 53 as a linear combination of 6731 and 5777 as follows, starting with the next to the last equation.

$$\begin{aligned} 53 &= 5777 - 6 \cdot 954 \\ &= 5777 - 6 \cdot (6731 - 1 \cdot 5777) = 7 \cdot 5777 - 6 \cdot 6731 \end{aligned}$$

Thus, the GCD of 6731 and 5777 is a linear combination of them.

Here's the general situation to find $\text{GCD}(m, n)$ as a linear combination of m and n . Let's suppose that $m > n$ to begin with. We divide n into m and get a quotient of q_1 and remainder of r_1 , that is

$$m = q_1 n + r_1,$$

with r_1 between 1 and n . Then we work with n and r_1 instead of m and n . Divide r_1 into n to get a quotient of q_2 and a remainder of r_2 , that is,

$$n = q_2 r_1 + r_2.$$

And we keep going until eventually we get a remainder of 0.

$$\begin{aligned} r_1 &= q_3 r_2 + r_3 \\ r_2 &= q_4 r_3 + r_4 \\ &\vdots \\ r_{s-3} &= q_{s-1} r_{s-2} + r_{s-1} \\ r_{s-2} &= q_s r_{s-1} + 0 \end{aligned}$$

We have

$$m > n > r_1 > r_2 > \cdots > r_{s-1}$$

and r_{s-1} is d , the GCD we're looking for.

Each equation finds a remainder as a linear combination of the previous two remainders. Starting with the next to the last equation, we can find $d = r_{s-1}$ as a linear combination of r_{s-2} and r_{s-3} . The equation before that gives r_{s-2} in terms of r_{s-3} and r_{s-4} , so we can also get d in terms of r_{s-3} and r_{s-4} . Working our way back up, we can eventually get d as a linear combination of m and n .

Thus, we've shown the following theorem.

Theorem 1.44 (Extended Euclidean algorithm). The greatest common divisor $d = \text{GCD}(m, n)$ of m and n is a linear combination of m and n . That is, there exist integers a and b such that

$$d = am + bn.$$

Now that we have the major theorems on GCDs, there are a few more fairly elementary properties of GCDs that are straightforward to prove, such as these.

Theorem 1.45.

$$\text{GCD}(a, b + ka) = \text{GCD}(a, b).$$

$$\text{GCD}(ak, bk) = k\text{GCD}(a, b).$$

$$\text{If } d = \text{GCD}(a, b) \text{ then } \text{GCD}(a/d, b/d) = 1.$$

Exercise 14. Prove the statements in the theorem.

Greatest common divisors of more than two numbers The GCD of more than two numbers is defined the same way as for two numbers: the GCD of a set of numbers the largest number that divides them all. For example, $\text{GCD}(14, 49, 91) = 7$. To find a GCD of three numbers, a , b , and c , first find $d = \text{GCD}(a, b)$, then find $e = \text{GCD}(d, c)$. Thus,

$$\text{GCD}(a, b, c) = \text{GCD}(\text{GCD}(a, b), c),$$

a statement that is easy to show.

Pairwise relatively prime numbers A set of numbers is said to be *pairwise relatively prime* or *pairwise coprime* if any two of them are relatively prime. For instance, 15, 22, and 49 are three pairwise relatively prime numbers. Thus, a , b , and c are pairwise relatively prime when

$$\text{GCD}(a, b) = \text{GCD}(a, c) = \text{GCD}(b, c) = 1.$$

Note that $\text{GCD}(a, b, c)$ can be 1 without a , b , and c being pairwise relatively prime. For instance, $\text{GCD}(6, 10, 15) = 1$, but $\text{GCD}(6, 10) = 2$, $\text{GCD}(6, 15) = 3$, and $\text{GCD}(10, 15) = 5$.

Least common multiples The *least common multiple* of a set of positive integers is the smallest positive integer that they all divide. It is easy to show that the greatest common divisor of two integers times their least common multiple equals their product.

$$\text{GCD}(a, b) \text{ LCM}(a, b) = ab.$$

Least common multiples can be used to sum common fractions. For example, to add $\frac{5}{6} + \frac{4}{15}$, note that the least common multiple of 6 and 15 is 30, so each fraction can be expressed with the least common denominator 30 as $\frac{25}{30} + \frac{8}{30} = \frac{25+8}{30} = \frac{33}{30}$. Even using least common denominators, it may be that the sum can be simplified as it can in this case to $\frac{11}{10}$.

1.5 The fundamental theorem of arithmetic

We proved above that every natural number could be factored as a product of primes. But we want more than existence, we want uniqueness. We need to prove that there is only one way that it can be factored as a product of primes.

The unique factorization theorem, a.k.a., the fundamental theorem of arithmetic.

Now, in order to make this general statement valid we have to extend a little bit what we mean by a product. For example, how do you write a prime number like 7 as a product of primes? It has to be written as the product 7 of only one prime. So we will have to accept a single number as being a product of one factor.

Even worse, what about 1? There are no primes that divide 1. One solution is to accept a product of no factors as being equal to 1. It's actually a reasonable solution to define the empty product to be 1, but until we find another need for an empty product, let's wait on that and restrict this unique factorization theorem to numbers greater than 1. So, here's the statement of the theorem we want to prove.

Theorem 1.46 (Unique factorization theorem). Each integer n greater than 1 can be uniquely factored as a product of primes. That is, if n equals the product $p_1 p_2 \cdots p_r$ of r primes, and it also equals the product $q_1 q_2 \cdots q_s$ of s primes, then the number of factors in the two products is the same, that is $r = s$, and the two lists of primes p_1, p_2, \dots, p_r and q_1, q_2, \dots, q_s are the same apart from the order the listings.

We'll prove this by using the strong form of mathematical induction. The form that we'll use is this:

In order to prove a statement $S(n)$ is true for all numbers, prove that $S(n)$ follows from the assumption that $S(k)$ is true for all $k < n$.

This principle of induction appears to be stronger than the one we've used before, but, in fact, it is equivalent to it. It's really the same as the minimization principle (i.e. well-ordering principle) applied to the negation of the statement. The advantage in using it is that a proof by contradiction is not needed making the proof more understandable.

Proof. We'll prove the unique factorization theorem in two cases. Case 1 will be where n is a prime number itself. Case 2 will be where n is composite.

Case 1: Suppose that n is a prime number. The only way that a prime number can be written as a product of primes is as itself; otherwise it would not be prime, but composite.

Case 2: Suppose that n is a composite number equal to both products of primes $p_1 p_2 \cdots p_r$ and $q_1 q_2 \cdots q_s$. Note that since n is composite, both r and s are at least 2; otherwise it would not be composite, but prime.

Now look at one of the primes, say p_1 . It divides n , so it divides the product of the other primes $q_1 q_2 \cdots q_s$. We suspect that that implies it has to be one of those other primes. Let's put that off for a bit; that is, logically before we prove this theorem, we need to prove another theorem, listed next, that if a prime divides a product of primes, then it is one of those primes; but we'll actually do that next. Assuming we've done that, then we can conclude that p_1 is one of the q_i 's. We can reorder the product $q_1 q_2 \cdots q_s$ to make it so that q_1 equals p_1 . Now, since $p_1 p_2 \cdots p_r = q_1 q_2 \cdots q_s$ and the first factors of the two products are equal, therefore $p_2 \cdots p_r = q_2 \cdots q_s$. Now, by our new induction principle, these are two prime factorizations of a number smaller than n , and hence are the same, except for their order. Therefore, they have the same number of factors, that is, $r = s$, and all the factors are the same except for their order. And the number n is that product times p_1 , which equals q_1 , therefore the original two products, $p_1 p_2 \cdots p_r$ and $q_1 q_2 \cdots q_s$, are the same except for order. Q.E.D.

Well, that finished the proof except we have to prove another theorem first, namely, the following one.

Theorem 1.47. If a prime divides a product of primes $q_1q_2 \dots q_s$, then it equals one of the primes q_1, q_2, \dots, q_s .

We could do that, but we we'll prove a slightly stronger theorem, namely, the following one.

Theorem 1.48. If a prime divides a product of numbers $b_1b_2 \dots b_s$, then it divides one of the numbers b_1, b_2, \dots, b_s .

Now the reason this theorem implies the previous theorem is because if a prime p divides a product of primes $q_1q_2 \dots q_s$, then it divides one of the primes q_1, q_2, \dots, q_s , but the only way that one prime can divide another is if it equals the other.

Proof. A product of s numbers $b_1b_2 \dots b_s$ is actually a series of binary products. It's b_1 times $b_2 \dots b_s$, and $b_2 \dots b_s$ is b_2 times $b_3 \dots b_s$, etc, where the last product is $b_{s-1}b_s$ is the product of b_{s-1} times b_s . That means that if we knew the following theorem, then, using ordinary induction, we could conclude this one. Q.E.D.

Theorem 1.49. If a prime divides a product of two numbers, then it divides one of the numbers.

Now, we could prove this theorem directly, but it turns out that there is a slightly stronger version that we can use in other places, so let's prove it, the one listed next, instead, and show this theorem follows from it.

Theorem 1.50. If n and a are relatively prime, and $n|ab$, then $n|b$.

Proof that this theorem implies implies the previous one. Suppose that a prime p divides ab . If p doesn't divide a , then it's relatively prime to a , so by this theorem, it divides b . Therefore, either $p|a$ or $p|b$. Q.E.D.

Proof of this theorem. Suppose that $\text{GCD}(n, a) = 1$. Then, by the extended Euclidean algorithm, 1 is a linear combination of n and a , that is, there exist integers t and u such that

$$1 = tn + ua.$$

Multiply that equation by b to get

$$b = tnb + uab.$$

Now, if $n|ab$, then n divides the right hand side of the equation, but that equals the left hand side, so $n|b$. Q.E.D.

Comment 1.51. Typically in a mathematics book those theorems that come first logically are presented first. Here we started with our goal and discovered the theorems that were needed to prove the goal. (Actually, I made the list longer than it needed to be by strengthening a couple of them because the stronger versions are more useful, something you can only tell with hindsight.)

The advantage to presenting theorems in their logical order is that it is easier to follow the logic. The disadvantage is that the motivation for the preliminary theorems is not apparent until the final theorem, the interesting one, is reached.

Usually when we write the prime factorization of a number, we'll use exponents on those primes that are repeated. For instance, the number 40 had the prime factorization $2 \cdot 2 \cdot 2 \cdot 5$. An abbreviated form for this factorization is $2^3 \cdot 5$. We say that the prime 2 occurs with multiplicity 3, while the prime 5 occurs with multiplicity 1. The multiplicities are the exponents. So, in general, a number n has the prime factorization

$$n = p_1^{e_1} p_2^{e_2} \cdots p_k^{e_k}$$

where the primes p_1, p_2, \dots, p_k are all distinct, and their multiplicities are the exponents e_1, e_2, \dots, e_k , respectively.

These exponents are called the orders of the primes in n . The *order* of p in n be the exponent of p in the prime factorization of n , denoted $\text{ord}_p a$.

Immediate corollaries to the unique factorization theorem. A corollary is a theorem that logically follows very simply from a theorem. Sometimes it follows from part of the proof of a theorem rather than from the statement of the theorem. In any case, it should be easy to see why it's true. We can draw a couple of corollaries from the unique factorization theorem.

Corollary 1.52. The only primes that can divide a number n are the ones that appear in its prime factorization $p_1^{e_1} p_2^{e_2} \cdots p_k^{e_k}$.

Corollary 1.53. If the prime factorizations of m and n are $m = p_1^{e_1} p_2^{e_2} \cdots p_k^{e_k}$ and $n = p_1^{f_1} p_2^{f_2} \cdots p_k^{f_k}$ (where here some of the e_i 's and f_i 's may equal 0 so we can use the same list of primes for both numbers), then their greatest common divisor $d = \text{GCD}(m, n)$ has the prime factorization $d = p_1^{g_1} p_2^{g_2} \cdots p_k^{g_k}$ where each exponent g_i is the minimum of the corresponding exponents e_i and f_i .

As an example of the last corollary, if $m = 1260 = 2^2 3^2 5^1 7^1$ and $n = 600 = 2^3 3^1 5^2$, then their GCD is $d = 2^2 3^1 5^1 = 60$.

1.6 Polynomials.

We'll frequently use polynomials in our study of fields and rings. We'll only consider polynomials with coefficients in fields and commutative rings, not with coefficients in noncommutative rings.

We won't formally define polynomials. For now, we'll only look at polynomials in one variable x , but later in section 3.10.4 we'll look at polynomials in two or more variables.

Informally a *polynomial* $f(x)$ with coefficients in a commutative ring R is an expression

$$f(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$$

where each coefficient $a_i \in R$. We'll assume that the leading coefficient a_n is not zero so that $\deg f$, the degree of the polynomial, is n . When a_n is zero, the polynomial is called a *monic* polynomial.

It's convenient to denote a polynomial either by f or by $f(x)$. If the variable x is referred to somewhere nearby, then I'll use $f(x)$, otherwise I'll just use f . For instance, if I want to multiply two polynomials f and g together, I'll write fg , but if I want to multiply f by $x^2 - 3x + 2$, I'll write $f(x)(x^2 - 3x + 2)$ or $f(x) \cdot (x^2 - 3x + 2)$.

A *root* of a polynomial is an element a of R such that $f(a) = 0$, that is, it's a solution of the polynomial equation $f(x) = 0$.

The set of all polynomials with coefficients in a commutative ring R is denoted $R[x]$. It has addition, subtraction, and multiplication, and satisfies the requirements of a ring, that is, it has addition, subtraction, and multiplication with the usual properties. $R[x]$ is called the *ring of polynomials with coefficients in R* . Note that $R[x]$ doesn't have reciprocals even when R is a field, since x has no inverse in $R[x]$. Therefore, $R[x]$ is not a field. Nonetheless, the ring R is a subring of the ring $R[x]$ since we can identify the constant polynomials as the elements of R .

1.6.1 Division for polynomials

Although $R[x]$ doesn't have reciprocals, it does have a division algorithm, at least when the divisor is a monic polynomial.

Theorem 1.54 (The division algorithm for polynomials over a ring). Let R be a commutative ring and $R[x]$ its polynomial ring in one variable. Let f be a polynomial (the dividend) and g a monic polynomial (the divisor). Then there exist unique polynomials q (the quotient) and r (the remainder) such that $f = qg + r$ where either $r = 0$ or $\deg r < \deg g$.

Proof of existence. One case is when $f = 0$ or $\deg f < \deg g$. Since the dividend already has a lower degree, the quotient $q = 0$ and the remainder $r = f$.

That leaves the case when $\deg f \geq \deg g$. We'll prove it by induction on $n = \deg f$ where the base case is $n = 0$. That's the case where f and g are both constants in the ring R , but g is monic, so $g = 1$. Then $q = f$ and $r = 0$.

Now for the inductive step. We'll assume the inductive hypothesis that the theorem is correct for all polynomials f of degree less than n and show it's true for those of degree n . Let

$$f(x) = a_0 + a_1 x + \cdots + a_n x^n \quad \text{and} \quad g(x) = b_0 + b_1 x + \cdots + b_{m-1} x^{m-1} + x^m \quad \text{where } n \geq m.$$

The polynomial $f_1(x) = f(x) - a_n x^{n-m} g(x)$ has a 0 coefficient for x^n , so its degree is less than n . By inductive hypothesis, there are polynomials q_1 and r_1 such that $f_1 = q_1 g + r_1$ where $r_1 = 0$ or $\deg r_1 < \deg g$. Equating the right sides of the two equations involving f_1 , we may conclude that

$$f(x) = (a_1(x) + a_n x^{n-m})g(x) + f_1(x).$$

That gives us the desired representation $f(x) = q(x)g(x) + r(x)$, finishing the inductive proof for the existence half of the proof. Q.E.D.

Proof of uniqueness. Suppose there are also polynomials q' and r' such that $f = q'g + r'$ where either $r' = 0$ or $\deg r' < \deg g$. We'll show $r = r'$ and $q = q'$.

Since $f = qg + r$ and $f = q'g + r'$, therefore $qg + r = q'g + r'$ so $r - r' = g(q' - q)$. Suppose that $r \neq r'$. Then $q' - q \neq 0$, and since g is a monic polynomial, therefore $\deg g(q - q') \geq \deg g$. Therefore $\deg(r - r') \geq \deg g$. But $\deg(r - r') < \deg g$ since both r and r' have degree less than $\deg g$, a contradiction. Therefore, $r = r'$.

Now we have $0 = g(q' - q)$, but g is monic, so $q' - q = 0$, and $q = q'$. Q.E.D.

If R happens to be a field, there is a stronger version of the theorem that doesn't require g to be a monic polynomial.

Theorem 1.55 (The division algorithm for polynomials over a field). Let F be a field and $F[x]$ its polynomial ring in one variable. Let f be a polynomial (the dividend) and g a nonzero polynomial (the divisor). Then there exist unique polynomials q (the quotient) and r (the remainder) such that $f = qg + r$ where either $r = 0$ or $\deg r < \deg g$.

Exercise 15. Prove the above theorem. Hint: divide g by its leading coefficient and use the division algorithm for polynomials over a ring. There will still be two parts, one for existence and one for uniqueness.

The remainder theorem and factor theorem. The remainder theorem is something that's frequently covered in high school algebra classes. It says when you divide a polynomial f by $x - a$, the remainder is $f(a)$. It works in general for polynomials with coefficients in an arbitrary ring.

Theorem 1.56 (Remainder theorem). Let R be a commutative ring and $R[x]$ its polynomial ring. For $f \in R[x]$ and $a \in R$, there is a polynomial q such that $f(x) = (x - a)q(x) + f(a)$.

Proof. Apply the division algorithm for $g(x) = x - a$. Then $f(x) = (x - a)q(x) + r$ where r is a constant. Setting x to a , we conclude $f(a) = r$. Q.E.D.

The factor theorem is a corollary of the remainder theorem. —

Theorem 1.57 (Factor theorem). For $f \in R[x]$ and $a \in R$, a is a root of f if and only if $(x - a)$ divides $f(x)$.

Further properties of polynomials. There are a couple more properties of polynomials that apply only when the ring is a field or an integral domain. As described later in section 3.1.3, an integral domain is a commutative ring in which $0 \neq 1$ that satisfies one of the two equivalent conditions: it has no zero-divisors, or it satisfies the cancellation law. Thus, fields are special cases of integral domains.

One property is that a polynomial of degree n has at most n roots.

Theorem 1.58. The number of roots of a nonzero polynomial with coefficients in an integral domain is at most the degree of the polynomial.

Proof. We'll prove this by induction on n , the degree of the polynomial f .

If $n = 0$, then f is a constant, but it's not the zero constant, so it has no roots.

Assume the inductive hypothesis, namely, the theorem holds for all functions of degree n . We'll show it holds for each function f of degree $n + 1$. If f has no roots, then the theorem is true, so let r be a root of f . By the factor theorem, $f(x) = (x - r)q(x)$, where the degree of the quotient q equals n .

We'll show every other root $r' \neq r$ of f is also a root of q . Since r' is a root, therefore $0 = f(r') = (r' - r)q(r)$. Now $r' - r$ is not 0, and the ring is an integral domain which has no zero-divisors, therefore $0 = q(r)$. Thus all other roots of f are roots of q .

Since $\deg g = n$, by the inductive hypothesis, g has at most n roots, therefore f has at most $n + 1$ roots.

That completes the proof by induction.

Q.E.D.

Exercise 16. An example of a ring that is not an integral domain is \mathbf{Z}_8 . Show that the quadratic polynomial $f(x) = x^2 - 1$ in $\mathbf{Z}_8[x]$ has more than two roots in \mathbf{Z}_8 .

A couple of corollaries for polynomials with coefficients in an integral domain follow from the previous theorems.

Corollary 1.59. If $\deg f = n$, and a_1, a_2, \dots, a_n are n distinct roots of f , then

$$f(x) = a(x - a_1)(x - a_2) \cdots (x - a_n)$$

where a is the leading coefficient of f .

Corollary 1.60. If two monic polynomials f and g both of degree n have the same value at n places, then they are equal.

1.6.2 Roots of unity and cyclotomic polynomials

Definition 1.61 (Root of unity). A *root of unity*, also called a root of 1 is a complex number such that when raised to some positive integer power yields 1. If $z^n = 1$, then z is called an n^{th} root of unity. If n is the smallest positive integer power such that $z^n = 1$, then n is called a n^{th} *primitive* root of unity.

Among the real numbers, the only roots of unity are 1 and -1 . 1 is the only first primitive root of unity and -1 is the only primitive second root of unity.

The n^{th} roots of unity are equally spaced around the unit circle separated by angles of $2\pi/n$. See figure 3.4 for the primitive seventh roots of unity on the unit circle.

An n root of unity z is a root of the polynomial $z^n - 1$, but not all roots of such a polynomial are primitive. For example, roots of the polynomial $z^2 - 1$ are second roots of unity, but 1, being one of those two roots, is not a primitive second root of unity.

Example 1.62 (Sixth roots of unity). The sixth roots of unity are roots of the polynomial $z^6 - 1$. This polynomial factors as $(z^3 + 1)(z^3 - 1) = (z^2 - z + 1)(z + 1)(z^2 + z + 1)(z - 1)$. Of course, two of the roots of this polynomial are 1 and -1 which account for the factors $x - 1$ and $x + 1$. The roots of the factor $z^2 - z - 1$ are also roots of $z^3 - 1$, so are cube roots of unity, in fact, they're the two primitive third roots of unity. Those roots are $z = \frac{1}{2}(-1 + i\sqrt{3})$. If you call one of them $\omega = \frac{1}{2}(-1 + i\sqrt{3})$, then the other one is $\omega^2 = \frac{1}{2}(-1 - i\sqrt{3})$. You can see

them displayed in the complex plane in figure 3.3 which illustrates the lattice of Eisenstein integers.

The roots of the other factor $z^2 + z + 1$ are $z = \frac{1}{2}(1 + i\sqrt{3})$. They are the two primitive sixth roots of unity. Notice that they are $\omega + 1$ and $\omega^2 + 1$.

So, altogether, there are six sixth roots of unity. Two are primitive sixth roots, two are primitive third roots, one is a primitive second root, and one is a primitive first root.

Among the five fifth roots of unity, one of them, $z = 1$, is not primitive, the other four are. They are roots of the polynomial $\Phi_5(z) = \frac{z^5 - 1}{z - 1} = z^4 + z^3 + z^2 + z + 1$.

If z is a primitive n^{th} root of unity, then the entire list of n^{th} roots is $1, z, z^2, \dots, z^{n-1}$. The root z^k won't be primitive if there is a common divisor of n and k . That leaves only $\varphi(n)$ of the roots to be primitive, where $\varphi(n)$ is the number of positive integers less than n that are relatively prime to n . See definition 1.39 for a definition of Euler's totient function φ .

Definition 1.63 (Cyclotomic polynomial). The polynomial $\Phi_n(z) = \prod_{k=1}^{\varphi(n)} (z - z_k)$, where $z_1, z_2, \dots, z_{\varphi(n)}$ are the primitive n^{th} roots of unity, is called the n^{th} cyclotomic polynomial.

There are two primitive third roots of unity as mentioned in the example above, so $\Phi_3(z) = z^2 - z + 1$. There are also two primitive sixth roots, and $\Phi_6(z) = z^2 + z - 1$.

When p is a prime number, then $\Phi(p)$ has degree $\varphi(p) = p - 1$. Its value is $\Phi(p) = \frac{z^p - 1}{z - 1} = z^{p-1} + \dots + z + 1$.

Here's a short table of the first few cyclotomic polynomials.

n	$\Phi(n)$	n	$\Phi(n)$
1	$z - 1$	9	$z^6 + z^3 + 1$
2	$z + 1$	10	$z^4 - z^3 + z^2 - z + 1$
3	$z^2 + z + 1$	11	$z^{10} + z^9 + \dots + z + 1$
4	$z^2 + 1$	12	$z^4 - z^2 + 1$
5	$z^4 + z^3 + z^2 + z + 1$	13	$z^{12} + z^9 + \dots + z + 1$
6	$z^2 - z + 1$	14	$z^6 - z^5 + z^4 - z^3 + z^2 - z + 1$
7	$z^6 + z^5 + z^4 + z^3 + z^2 + z + 1$	15	$z^8 - z^7 + z^5 - z^4 + z^3 - z + 1$
8	$z^4 + 1$	16	$z^8 + 1$

It's interesting that the only coefficients that appear in the first one hundred cyclotomic polynomials are 0, 1, and -1 .

We'll use cyclotomic polynomials in section 3.10.3.

Chapter 2

Fields

Informally, a field is a set equipped with four operations—addition, subtraction, multiplication, and division that have the usual properties.

We'll study rings in chapter 3, which are like fields but need not have division.

2.1 Introduction to fields

A *field* is a set equipped with two binary operations, one called *addition* and the other called *multiplication*, denoted in the usual manner, which are both commutative and associative, both have identity elements (the additive identity denoted 0 and the multiplicative identity denoted 1), addition has inverse elements (the inverse of x being denoted $-x$), multiplication has inverses of nonzero elements (the inverse of x being denoted $\frac{1}{x}$), multiplication distributes over addition, and $0 \neq 1$.

Three fields that you already know are the field of real numbers \mathbf{R} , the field of rational numbers \mathbf{Q} , and the field of complex numbers \mathbf{C} .

We'll see that there are many other fields. When we have a generic field, will use a capital F to denote it.

2.1.1 Definition of fields

Here's a more complete definition.

Definition 2.1 (field). A *field* F consists of

1. a set, also denoted F and called the *underlying set* of the field;
2. a binary operation $+$: $F \times F \rightarrow F$ called *addition*, which maps an ordered pair $(x, y) \in F \times F$ to its *sum* denoted $x + y$;
3. another binary operation \cdot : $F \times F \rightarrow F$ called *multiplication*, which maps an ordered pair $(x, y) \in F \times F$ to its *product* denoted $x \cdot y$, or more simply just xy ;
such that
4. addition is commutative, that is, for all elements x and y , $x + y = y + x$;

5. multiplication is commutative, that is, for all elements x and y , $xy = yx$;
6. addition is associative, that is, for all elements x , y , and z , $(x + y) + z = x + (y + z)$;
7. multiplication is associative, that is, for all elements x , y , and z , $(xy)z = x(yz)$;
8. there is an additive identity, an element of F denoted 0 , such that for all elements x , $0 + x = x$;
9. there is a multiplicative identity, an element of F denoted 1 , such that for all elements x , $1x = x$;
10. there are additive inverses, that is, for each element x , there exists an element y such that $x + y = 0$; such a y is called the *negation* of x ;
11. there are multiplicative inverses of nonzero elements, that is, for each nonzero element x , there exists an element y such that $xy = 1$; such a y is called a *reciprocal* of x ;
12. multiplication distributes over addition, that is, for all elements x , y , and z , $x(y + z) = xy + xz$; and
13. $0 \neq 1$.

The conditions for a field are often call the *field axioms*.

Caveat: We're using the terminology and notation of arithmetic that we use for numbers, but the elements of our fields need not be numbers; often they will be, but sometimes they won't.

Note that we'll use the standard notational conventions on precedence for all fields so we don't have to fully parenthesize every expression. Multiplication and division have a higher precedence than addition and subtraction, so that, for example, $x - y/z$ means $x - (y/z)$, not $(x - y)/z$. Also, operations are executed from left to right, so that $x - y - z$ means $(x - y) - z$, not $x - (y - z)$. (Usually operations are executed from left to right, but an exception is that exponentiation is executed from right to left, so that x^{m^n} means $x^{(m^n)}$, not $(x^m)^n$.)

Commutativity and associativity of addition imply that terms can be added in any order, so of course we won't put parentheses when we're adding more than two terms together. Likewise for multiplication.

Although in parts 10 and 11 of the definition only the existence of an additive and multiplicative inverses is required, you can easily show uniqueness follows from the definition. Once that is done we can note that the additive inverse of x is called *the negation* of x and denoted $-x$, and the multiplicative inverse of x , when x is not 0, is called *the reciprocal* of x and denoted $1/x$, $\frac{1}{x}$, or x^{-1} .

2.1.2 Subtraction, division, multiples, and powers

With the help of negation, we can define subtraction as follows. The *difference* of two elements x and y is defined as $x - y = x + (-y)$.

Likewise, with the help of reciprocation, we can define division. The *quotient* of an element x and a nonzero element y is xy^{-1} , denoted x/y or $\frac{x}{y}$. The expected properties of subtraction

and division all follow from the definition of fields. For instance, multiplication distributes over subtraction, and division by z distributes over addition and subtraction.

Likewise, we can define integral multiples of elements in a field. First, we'll define nonnegative multiples inductively. For the base case, define $0x$ as 0 . Then define $(n+1)x$ as $x+nx$ when n is a nonnegative integer. Thus nx is the sum of n x 's. For instance, $3x = x + x + x$. Then if $-n$ is a negative integer, we can define $-nx$ as $-(nx)$. The usual properties of multiples, like $(m+n)x = mx + nx$ will, of course, hold.

Furthermore, we can define integral powers of x . Define x^1 as x for a base case, and inductively for nonnegative n , define x^{n+1} as xx^n . Thus nx is the product of n x 's. For instance, $x^3 = xxx$. Next, define x^0 as 1 , so long as $x \neq 0$. (0^0 should remain undefined, but for some purposes, especially in algebra, it's useful to define 0^0 to be 1 .) Finally, if $-n$ is positive and $x \neq 0$, define x^{-n} as $(x^n)^{-1}$. The usual properties of integral powers hold, like $x^{m+n} = x^m x^n$ and $(xy)^n = x^n y^n$.

2.1.3 Properties that follow from the axioms

There are numerous useful properties that are logical consequences of the axioms. Generally speaking, the list of axioms should be short, if not minimal, and any properties that can be proved should be proved. Here's a list of several things that can be proved from the axioms. We'll prove a few in class, you'll prove some as homework, and we'll leave the rest. (They make good questions for quizzes and tests.)

In the following statements, unquantified statements are meant to be universal with the exception that whenever a variable appears in a denominator, that variable is not to be 0 .

Exercise 17. Prove that 0 is unique. That is, there is only one element x of a field that has the property that for all y , $x + y = y$. (The proof that 1 is unique is similar.)

Exercise 18. Prove that each number has only one negation. That is, for each x there is only one y such that $x + y = 0$. (The proof that reciprocals of nonzero elements are unique is similar.)

Exercise 19. Prove that the inverses of the identity elements are themselves, that is, $-0 = 0$, and $1^{-1} = 1$.

Exercise 20. Prove that multiplication distributes over subtraction: $x(y - z) = xy - xz$.

Exercise 21. Prove that 0 times any element in a field is 0 : $0x = 0$.

Exercise 22. Prove the following properties concerning multiplication by negatives: $(-1)x = -x$, $-(-x) = x$, $(-x)y = -(xy) = x(-y)$, and $(-x)(-y) = xy$.

Exercise 23. Prove the following properties concerning reciprocals: $(x^{-1})^{-1} = x$, and $(xy)^{-1} = x^{-1}y^{-1}$.

Exercise 24. Prove that when y and z are both nonzero that $\frac{x}{y} = \frac{w}{z}$ if and only if $xz = yw$.

Exercise 25. Prove the following properties concerning division:

$$\begin{aligned} \frac{x}{y} \pm \frac{w}{z} &= \frac{xz \pm yw}{yz}, \\ \frac{x}{y} \frac{w}{z} &= \frac{xw}{yz}, \text{ and} \\ \frac{x}{y} \bigg/ \frac{w}{z} &= \frac{xz}{yw}. \end{aligned}$$

Assume that any time a term appears in a denominator that it does not equal 0.

Exercise 26. Prove that if $xy = 0$, then either $x = 0$ or $y = 0$.

2.1.4 Subfields

Frequently we'll find one field contained in another field. For instance, the field of rational numbers \mathbf{Q} is part of the field of real numbers \mathbf{R} , and \mathbf{R} is part of the field of complex numbers \mathbf{C} . They're not just subsets, $\mathbf{Q} \subset \mathbf{R} \subset \mathbf{C}$, but they have the same operations. Here's the precise definition of subfield.

Definition 2.2 (subfield). A field E is a *subfield* of a field F if

1. the underlying set of E is a subset of the underlying set of F ;
2. the addition operation $+_E$ on E is the restriction of the addition operation $+_F$ on F , that is, for all x and y in E , $x +_E y = x +_F y$; and
3. the multiplication operation \cdot_E on E is the restriction of the multiplication operation \cdot_F on F , that is, for all x and y in E , $x \cdot_E y = x \cdot_F y$.

When E is a subfield of F , we'll also say that F is an *extension* of E .

When you know one field is a subfield of another, there's no need to subscript the operations since they are the same.

There is an alternate characterization of subfield. The proof of the following theorem is straightforward, but there are many steps.

Theorem 2.3. If a subset E of a field F has 0, 1, and is closed under addition, multiplication, negation, and reciprocation of nonzero elements, then E is a subfield of F .

The field of rational numbers \mathbf{Q} . When we're trying to find the smallest example of a field, it looks like it will have to be \mathbf{Q} . Later in section 2.2 we'll see that it's not the smallest! But here's an argument (which must have a flaw in it) which says we need all the rational numbers to be in any field F .

To begin with, 0 and 1 have to be in F . But we also have to have $1 + 1$ in F and we'll denote that 2, of course. And we'll need $1 + 1 + 1 = 2 + 1$ which we'll denote 3. And so forth, so we've got 0 and all the positive integers in F . We also need negations of them, so all the negative integers are in F , too. But a rational number m/n is just an integer m divided by a positive integer n , so we'll have to have all rational numbers in F . That shows that \mathbf{Q} is a subfield of F .

Thus, it looks like every field F includes the smallest field \mathbf{Q} , the field of rational numbers.

There's one minor flaw in the argument above, but let's not pick it apart right now. Pretty soon we'll look at fields that don't contain \mathbf{Q} .

2.1.5 Fields of rational functions

A *rational function* with coefficients in F is a quotient of two polynomials $\frac{f(x)}{g(x)}$. Rational functions do form a field, the field $F(x)$ of rational functions with coefficients in F . Notice that the ring of polynomials $F[x]$ is denoted with square brackets while the field of rational functions $F(x)$ is denoted with round parentheses.

For example, one rational function in $\mathbf{Q}(x)$ is $\frac{5x^2 - 3x + 1/2}{x^3 + 27}$.

Note that the field F is a subfield of the $F(x)$. Again, we can identify the constant rational functions as the elements of F . For example, \mathbf{Q} is a subfield of $\mathbf{Q}(x)$, and both \mathbf{R} and $\mathbf{Q}(x)$ are subfields of $\mathbf{R}(x)$.

Also, the the ring of polynomials with coefficients is a subring of the field of rational functions. That is $F \subseteq F[x] \subseteq F(x)$.

2.1.6 Vector spaces over arbitrary fields

When you studied vector spaces, you may have studied only vector spaces over the real numbers, although vector spaces over other fields might have been mentioned. In fact, vector spaces over an arbitrary field F have the same basic properties as vector spaces over \mathbf{R} .

The n -dimensional standard vector space F^n is defined the same way as \mathbf{R}^n except the n -tuples have coordinates in F . Addition and scalar multiplication are defined the same way for F^n as they are for \mathbf{R}^n .

Furthermore, matrices in $M_{m \times n}(F)$ are defined the same way as matrices in $M_{m \times n}(\mathbf{R})$ except the entries are in F instead of \mathbf{R} . The matrix operations are the same. You can use the same methods of elimination to solve a system of linear equations with coefficients in F or find the inverse of a matrix in $M_{n \times n}(F)$ if its determinant is nonzero. Determinants have the same properties. You can use methods of linear algebra to study geometry in F^n just as you can for \mathbf{R}^n (although it may not be possible to visualize what F^n is supposed to look like, and things like areas of triangles have values in F).

The abstract theory of finite dimensional vector spaces over F is the same, too. Linear independence, span, basis, dimension are all the same. Rank and nullity of a matrix are the same. Change of basis is the same.

Eigenvalues, eigenvectors, and eigenspaces may have problems over some fields. In fact, when you studied transformations $\mathbf{R}^n \rightarrow \mathbf{R}^n$, sometimes you had complex eigenvalues, and their only eigenvectors were in \mathbf{C}^n . Likewise when looking at transformations $F^n \rightarrow F^n$ and the eigenvalues aren't in F , you'll may have to go to some field extension F' of F to find them and to F'^n to find the eigenvectors.

Likewise, canonical forms for matrices will depend on F .

2.2 Cyclic rings and finite fields

In this section we'll look at fields that are finite, and we'll discover that \mathbf{Q} actually isn't the smallest field. Although they're smaller fields—they're finite—they won't be subfields of \mathbf{Q} .

First we'll look a bit at the concept of congruence modulo n , where n is a positive integer. Then we'll look at the ring of integers modulo n , denoted $\mathbf{Z}/n\mathbf{Z}$ or more simply \mathbf{Z}_n . We'll

see why they're called cyclic rings. Finally, we'll look at the case where n is prime, and we'll denote it p then, where \mathbf{Z}_p turns out to be a field, and we'll examine some of the cyclic fields.

Definition 2.4 (Congruence modulo n). Fix n , a positive integer. We say that two integers x and y are *congruent modulo n* if n evenly divides the difference $x - y$. We'll use the standard notation from number theory

$$x \equiv y \pmod{n}$$

to indicate that x is congruent to y modulo n , and the notation $n|m$ to indicate that the integer n divides the integer m (with no remainder). Then

$$x \equiv y \pmod{n} \quad \text{iff} \quad n|(x - y).$$

When n doesn't divide the difference $x - y$, we say a is not congruent to b , denoted $x \not\equiv y \pmod{n}$.

You're familiar with congruence modulo 12; it's what 12-hour clocks use.

The general theory of equivalence relations in section A.2.3.

Theorem 2.5. Congruence modulo n is an equivalence relation.

Proof. For reflexivity, $x \equiv x \pmod{n}$ holds since $n|(x - x)$.

For symmetry, we need to show that $x \equiv y \pmod{n}$ implies $y \equiv x \pmod{n}$. But if $n|(x - y)$, then $n|(y - x)$.

For transitivity, suppose that $x \equiv y \pmod{n}$ and $y \equiv z \pmod{n}$. Then $n|(x - y)$ and $n|(y - z)$, so there exist k and m such that $nk = x - y$ and $nm = y - z$. Therefore $n(k + m) = x - z$, showing that $n|(x - z)$. Hence $x \equiv z \pmod{n}$. Q.E.D.

2.2.1 The cyclic ring \mathbf{Z}_n

Definition 2.6 (Integers modulo n). The integers modulo n , \mathbf{Z}_n is the set of equivalence classes of integers under the equivalence relation which is congruence modulo n .

We'll denote these equivalence classes with square brackets subscripted by n . Thus, for instance, the element 0 in \mathbf{Z}_6 is really $[0]_6$, which we'll denote $[0]$ when modulo 6 is understood. This equivalence class is the set of all x such that $x \equiv 0 \pmod{6}$. This $[0]_6 = \{\dots, -18, -12, -6, 0, 6, 12, 18, \dots\}$. Likewise the element 1 in \mathbf{Z}_6 is really the equivalence class of 1, which is the set

$$[1]_6 = \{x \in \mathbf{Z} \mid x \equiv 1 \pmod{6}\} = \{\dots, -17, -11, -5, 1, 7, 13, 19, \dots\}.$$

Note that $[1]_6 = [7]_6 = [13]_6$ all name the same equivalence class.

An equation in equivalence classes, such as $[x]_6 + [3]_6 = [5]_6$, is the same thing as an congruence, $x + 3 \equiv 5 \pmod{6}$. The congruence notation is usually more convenient. When the modulus n is known by context, we'll dispense with the subscript n , and abusing notation, we'll frequently drop the square brackets.

There are two ways you can think about integers modulo n . One is to think of them as regular integers from 0 through $n - 1$, do the arithmetic modulo n , and adjust your answer so it's in the same range. For example, we can take $\mathbf{Z}_6 = \{0, 1, 2, 3, 4, 5\}$. Then, to do some

computation, say $5(2 - 4) \pmod{6}$, first compute $5(2 - 4)$ as integers to get -10 , and then, since $-10 \equiv 2 \pmod{6}$, say the answer is 2. That works very well for computation, but it's pretty messy when you're trying to do anything with variables or trying to prove anything in general.

A better way is to say that an element of \mathbf{Z}_n is named by an integer, but two integers name the same element of \mathbf{Z}_n if they're congruent modulo n . Thus, x and y name the same element of \mathbf{Z}_n if $x \equiv y \pmod{n}$. This will work because congruence modulo n is an equivalence relation as we saw earlier.

In any case, it helps conceptually to think of the elements of \mathbf{Z}_n as being arranged on a circle like we imagine the elements of \mathbf{Z} being arranged on a line. See figure 2.1 of a couple of cyclic rings \mathbf{Z}_n to see where the word "ring" came from.

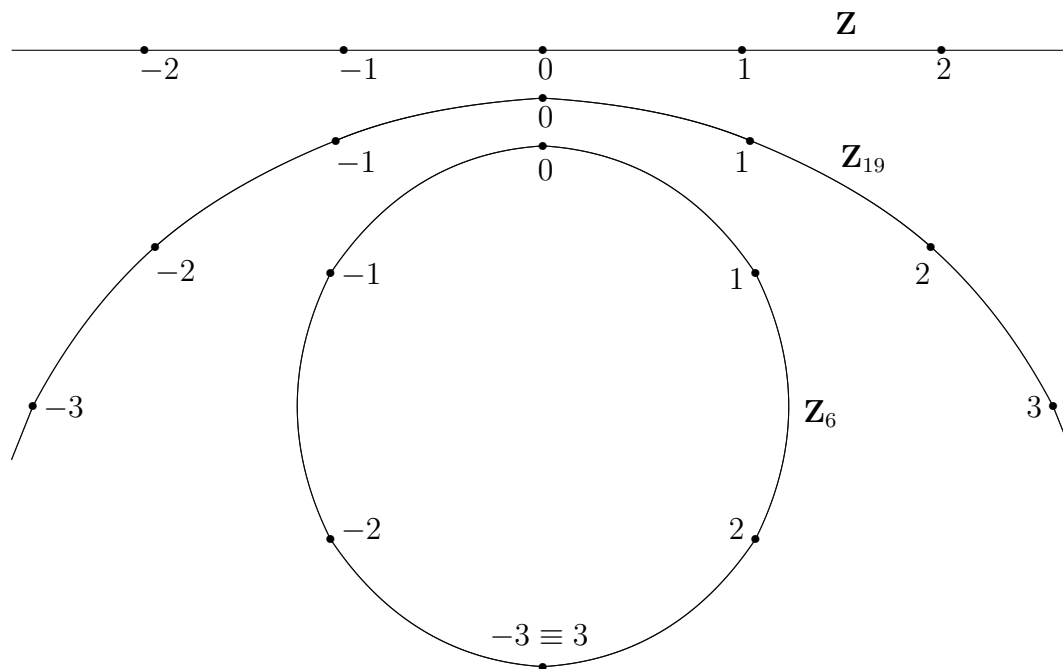


Figure 2.1: Cyclic rings $\mathbf{Z}_6, \mathbf{Z}_{19}, \mathbf{Z}$

The operations on \mathbf{Z}_n . Our equivalence relation is congruence modulo n , so our equivalence classes are also called congruence classes.

Congruence modulo n is more than just an equivalence relation; it works well with addition, subtraction, and multiplication, as you can easily show.

Theorem 2.7. If $x \equiv y \pmod{n}$, and $u \equiv v \pmod{n}$, then $x + u \equiv y + v \pmod{n}$, $x - u \equiv y - v \pmod{n}$, and $xu \equiv yv \pmod{n}$.

These properties will allow us to define a ring structure on \mathbf{Z}_n , as done below.

But congruence modulo n doesn't work so well with division. Although $6 \equiv 0 \pmod{6}$, it is not the case that $6/2 \equiv 0/2 \pmod{6}$. Thus, we can't expect that \mathbf{Z}_n will be a field, at least when $n = 6$.

Our job is to define addition, subtraction, and multiplication on \mathbf{Z}_n . Whenever a set is defined as a quotient set, that is, on equivalence classes, as \mathbf{Z}_n is, an extra step is required when defining an operation on it, as we'll see.

We would like to define addition on \mathbf{Z}_n by saying $[x] + [u] = [x + u]$, that is, the sum of the equivalence class of x and the equivalence class of u should be the equivalence class of $x + u$. But what if we named the equivalence class x by some other integer, say y , and the equivalence class of u by some other integer v ? How do we know that $[y + v]$ is the same equivalence class as $[x + u]$? We can state this question in a couple of other ways. How do we know

$$[x] = [y] \text{ and } [u] = [v] \text{ implies } [x + u] = [y + v]?$$

That asks the question: how do we know

$$x \equiv y \pmod{n} \text{ and } u \equiv v \pmod{n} \text{ implies } x + u \equiv y + v \pmod{n}?$$

That's one of the properties of congruence mentioned above. That property says addition on \mathbf{Z}_n is "well-defined".

Likewise, since multiplication works well with congruence,

$$x \equiv y \pmod{n} \text{ and } u \equiv v \pmod{n} \text{ imply } xu \equiv yv \pmod{n},$$

we can define multiplication on \mathbf{Z}_n by $[x] \cdot [u] = [xu]$.

Furthermore, all the ring axioms will be satisfied in \mathbf{Z}_n since they're satisfied in \mathbf{Z} . Thus, \mathbf{Z}_n is a ring, and it's called a *cyclic ring*.

The projection $\gamma : \mathbf{Z} \rightarrow \mathbf{Z}_n$. The function $\gamma : \mathbf{Z} \rightarrow \mathbf{Z}_n$ defined by $\gamma(k) = [k]$ maps an integer to its equivalence class modulo n . We defined addition and multiplication in \mathbf{Z}_n

$$[x + u] = [x] + [u] \quad \text{and} \quad [xu] = [x][u]$$

so γ preserves addition and multiplication. Furthermore, since $\gamma(1) = [1]$, it preserves 1. Therefore γ is a ring homomorphism. It is, of course, onto, so it is a ring epimorphism. It's called a *projection* or a *canonical homomorphism* to the quotient ring.

In section 3.6, we'll generalize this construction to rings besides \mathbf{Z} and their quotients, and we'll have projections for the generalizations, too.

The characteristic of a ring. What's weird about these cyclic rings is that if you start with 1 and add 1 over and over, you'll reach zero. For instance, in \mathbf{Z}_5 , we have $1+1+1+1+1 = 5 \equiv 0 \pmod{5}$. This corresponds to the geometric interpretation of these cyclic rings being shaped like rings.

Definition 2.8. If some multiple of 1 equals 0 in a ring, then the *characteristic* of the ring is the smallest such multiple. If no multiple of 1 equals 0, then the characteristic is said to be 0.

We're primarily interested in characteristics when we're talking about fields, and we'll see soon that the characteristic of a field is either 0 or a prime number.

Example 2.9. The characteristic of \mathbf{Z}_5 is 5, and, in general, the characteristic of a finite cyclic ring \mathbf{Z}_n is n .

2.2.2 The cyclic prime fields \mathbf{Z}_p

Since division doesn't work well with congruence, we can't expect \mathbf{Z}_n to always have reciprocals, so we don't expect it to be a field. Let's first see when an element in \mathbf{Z}_n is a unit. The term *unit* in a ring refers to an element x of the ring that does have a reciprocal. 1 is always a unit in a ring, and every nonzero element in a field is a unit.

Theorem 2.10. An element k in \mathbf{Z}_n is a unit if and only if k is relatively prime to n .

Proof. First, suppose that k is a unit in \mathbf{Z}_n . That means there exists l such that $kl \equiv 1 \pmod{n}$. Then $n \mid (kl - 1)$, and hence n is relatively prime to k .

Second, suppose that k is relatively prime to n . Then, by the extended Euclidean algorithm, their greatest common divisor, 1, is a linear combination of k and n . Thus, there are integers x and y so that $1 = xk + yn$. Then $1 \equiv xk \pmod{n}$, and k does have a reciprocal, namely x , in \mathbf{Z}_n . Thus k is a unit in \mathbf{Z}_n . Q.E.D.

Recall from definition 1.39 that the totient function $\varphi(n)$ denotes the number of positive integers less than n that are relatively prime to n .

Corollary 2.11 (Units in \mathbf{Z}_n). The number of units in \mathbf{Z}_n is $\phi(n)$.

Theorem 2.12. The cyclic ring \mathbf{Z}_n is a field if and only if n is prime.

Proof. Part of this theorem is a direct corollary of the previous one. Suppose n is prime. Then every nonzero element of \mathbf{Z}_n is relatively prime to n . Therefore, \mathbf{Z}_n is a field.

Next we'll show that if n is composite, the ring is not a field. Let n be the product of two integers m and k , both greater than 1. Then neither m nor k can have a reciprocal in \mathbf{Z}_n . Why not? Suppose that m^{-1} did exist in \mathbf{Z}_n . Then

$$\begin{aligned}(m^{-1}m)k &\equiv 1k \equiv k \pmod{n} \\ m^{-1}(mk) &\equiv m^{-1}n \equiv 0 \pmod{n}\end{aligned}$$

But $k \not\equiv 0 \pmod{n}$, a contradiction. So m^{-1} doesn't exist. Therefore, \mathbf{Z}_n is not a field. Q.E.D.

Corollary 2.13. The characteristic of a field is 0 or a prime number.

Proof. We'll show that if the characteristic n is finite, it must be a prime number. Suppose $n = st$. Then $0 = n \cdot 1 = (st) \cdot 1 = (s \cdot 1)(t \cdot 1)$. Therefore, either $s \cdot 1 = 0$ or $t \cdot 1 = 0$. But n is the smallest positive integer such that $n \cdots 1 = 0$, so either $s = n$ or $t = n$. Therefore n is prime. Q.E.D.

This proof works as well in integral domains introduced in section 3.1.3. This theorem will be mentioned again at that time.

Example 2.14. \mathbf{Z}_2 . Note that there is only one nonzero element, namely 1, and it is its own inverse. The addition and multiplication tables for \mathbf{Z}_2 are particularly simple.

$$\begin{array}{c|cc} + & 0 & 1 \\ \hline 0 & 0 & 1 \\ 1 & 1 & 0 \end{array} \qquad \begin{array}{c|cc} \cdot & 0 & 1 \\ \hline 0 & 0 & 0 \\ 1 & 0 & 1 \end{array}$$

Note that subtraction is the same as addition in \mathbf{Z}_2 since $x - y \equiv x + y \pmod{2}$.

Example 2.15. \mathbf{Z}_3 . Here, there are two nonzero elements, namely 1 and 2, but, for symmetry's sake, we'll call the two nonzero elements 1 and -1 . Note that each of these two are their own inverses. The addition and multiplication tables are still pretty simple.

+	-1	0	1
-1	1	-1	0
0	-1	0	1
1	0	1	-1

·	-1	0	1
-1	1	0	-1
0	0	0	0
1	-1	0	1

Example 2.16. \mathbf{Z}_{13} . What are the reciprocals of the 12 nonzero elements? We can name the nonzero elements as $\pm 1, \pm 2, \dots, \pm 6$. You can verify that this table gives their inverses.

x	± 1	± 2	± 3	± 4	± 5	± 6
x^{-1}	± 1	∓ 6	∓ 4	∓ 3	∓ 5	∓ 2

For instance, the reciprocal of 2 is -6 since $2(-6) \equiv -12 \equiv 1 \pmod{13}$.

These fields, \mathbf{Z}_p where p is prime, are the finite prime fields. But there are other finite fields.

Example 2.17. A field of order 9. We'll make an extension of \mathbf{Z}_3 to get a field of order 9. Note that -1 is not a square modulo 3. We can append $\sqrt{-1}$ to \mathbf{Z}_3 to get a field algebraic over it in exactly the same way we got \mathbf{C} from \mathbf{R} . Let's use i as an abbreviation for $\sqrt{-1}$, as usual. Then

$$\mathbf{Z}_3(i) = \{x + yi \mid x, y \in \mathbf{Z}_3\}$$

Addition, subtraction, and multiplication give us no problems. We just have to check that nonzero elements have inverses. That's exactly as before.

$$\frac{1}{x + yi} = \frac{x - yi}{(x + yi)(x - yi)} = \frac{x - yi}{x^2 + y^2} = \frac{x}{x^2 + y^2} + \frac{-y}{x^2 + y^2}i$$

Thus, if $x + yi$ is not 0 in $\mathbf{Z}_3(i)$, that is, not both of x and y are congruent to 0 modulo 3, then $x^2 + y^2 \not\equiv 0 \pmod{3}$, and the expression on the right gives $(x + yi)^{-1}$. Note that the characteristic of this field is 3 since $1 + 1 + 1$ is 0 in this field.

Exercise 27. You can construct a field of order 25 from \mathbf{Z}_5 , but it has to be done somewhat differently because $\sqrt{-1}$ already exists in \mathbf{Z}_5 since $(\pm 2)^2 = 4 \equiv -1$ in \mathbf{Z}_5 . The squares of the nonzero elements in \mathbf{Z}_5 include $1 = (\pm 1)^2$ and $4 = (\pm 2)^2$, but 2 is not among the squares. Show that the ring $\mathbf{Z}_5[\sqrt{2}]$ is a field by finding an inverse of a nonzero element $x + y\sqrt{2}$ where x and y are elements of \mathbf{Z}_5 but not both are 0. Hint: $(x - y\sqrt{2})(x + y\sqrt{2}) = x^2 - 2y^2$ cannot be 0.

In fact, there are finite fields of order p^n for each power of a prime p . These are called the Galois fields $GF(p^n)$. Note that cyclic prime field are the simplest Galois fields; \mathbf{Z}_p is $GF(p)$. The example constructed $GF(3^2)$ and the exercise $GF(5^2)$.

The proof that a finite field of characteristic p has to have p^n elements follows from the theorems in section 4.9.2 on Abelian groups. It's only dependent on addition in the finite field, not on multiplication.

2.2.3 Characteristics of fields, and prime fields

The characteristic of a ring was defined above, so we already have the definition for the characteristic of a field.

Those fields that have characteristic 0 all have \mathbf{Q} as a subfield. The flawed proof we saw earlier included the mistaken assumption that all the elements $0, 1, 2, \dots$ were distinct, which, as we've seen with these finite fields, isn't always the case. But we can correct the flawed proof to validate the following theorem. First, a definition.

Definition 2.18. A *prime field* is a field that contains no proper subfield. Equivalently, every element in it is a multiple of 1.

Theorem 2.19. Each field F has exactly one of the prime fields as a subfield. It will have \mathbf{Z}_p when it has characteristic p , but it will have \mathbf{Q} if it has characteristic 0.

The Frobenius endomorphism. Exponentiation to the power p has an interesting property when a commutative ring R has prime characteristic p :

$$(x + y)^p = x^p + y^p$$

There are various ways to prove this. For instance, you can show that the binomial coefficient $\binom{p}{k}$ is divisible by p when $1 < k < p$. Since $\binom{p}{k} = \frac{p!}{k!(n-k)!}$ and p divides the numerator but not the denominator, therefore it divides $\binom{p}{k}$.

This function $\varphi : R \rightarrow R$ defined by $\varphi(x) = x^p$ also preserves 1 and multiplication: $1^p = 1$ and $(xy)^p = x^p y^p$. Therefore, it is a ring endomorphism, called the *Frobenius* endomorphism.

We're most interested in the endomorphism when the ring is a field F of characteristic p . It's not particularly interesting when F is the prime field \mathbf{Z}_p because it's just the identity function then. For other finite fields of characteristic p it will be an automorphism—it's a bijection since it's an injection on a finite set—and it's not the identity function for those fields.

Example 2.20. In the example above of the Galois field $GF(3^2) = \mathbf{Z}_3(i)$, the characteristic of the field is 3, so $\varphi(x + yi) = (x + yi)^3 = x^3 + (yi)^3 = x^3 - y^3i = x - yi$. On the subfield \mathbf{Z}_3 , φ is the identity, but not on all of $GF(3^2) = \mathbf{Z}_3(i)$, since $\varphi(i) = -i$.

Exercise 28. Determine the value of $\varphi(\sqrt{2})$ in $GF(5^2)$.

2.3 Field Extensions, algebraic fields, the complex numbers

A lot of fields are found by extending known fields. For instance, the field of complex numbers \mathbf{C} is extended from the field of real numbers \mathbf{R} and $GF(3^2)$ is extended from $\mathbf{Z}_3 = GF(3)$. We'll look at the general case of extending fields by adding square roots to known fields, the smallest kind of extension, called a *quadratic* extension.

2.3.1 Algebraic fields

We've looked at some quadratic extensions of fields. Now we'll look at algebraic extensions in more detail.

Definition 2.21 (Algebraic and transcendental numbers). An *algebraic* number is a number that is a root of a polynomial with rational coefficients. If the polynomial is monic, then the algebraic number is an *algebraic integer*. A real number or a complex number that is not algebraic is called a *transcendental* number.

For instance, $x = \sqrt{2}$ is an algebraic number since it is the root of the polynomial $x^2 - 2$; in fact, it's an algebraic integer. On the other hand, $x = \sqrt{1/2}$ is a root of the polynomial $2x^2 - 1$, so it's an algebraic number, but not an algebraic integer.

There are many real numbers used in analysis that are transcendental. In 1873 Charles Hermite (1822–1901) proved that the number e is transcendental. It follows that many related numbers are transcendental such as e^2 and \sqrt{e} .

Definition 2.22 (Algebraic and transcendental field extensions). More generally, if x satisfies a polynomial equation $f(x) = 0$ where the polynomial f has coefficients in a field F , then we say x is *algebraic* over F . A field extension F' of F , all of whose elements are algebraic over F is said to be an *algebraic* extension of F . Field extensions that are not algebraic are called *transcendental* extensions. An algebraic extension of \mathbf{Q} is also called an *algebraic number field*, or more simply a number field.

In 1882 Lindemann extended Hermite's result to show that e^a is transcendental for all nonzero algebraic numbers a . Thus $e^{\sqrt{2}}$ is transcendental. More importantly, Lindemann's theorem shows that $\pi = e^i$ is transcendental, for if it were algebraic, then $e^{\pi i} = -1$ would be transcendental, which it isn't.

Weierstrass proved an even more general theorem in 1885. If a_1, \dots, a_n are distinct nonzero algebraic numbers, then the numbers e^{a_1}, \dots, e^{a_n} are algebraically independent meaning each e^{a_i} is transcendental over the field $\mathbf{Q}(e^{a_1}, \dots, e^{\hat{a}_i}, e^{a_n})$. The hat over e^{a_i} means that is omitted from the list.

Example 2.23. We know that the square root of 2, $\sqrt{2}$ is not a rational number. The field $\mathbf{Q}(\sqrt{2})$ is the smallest field that contains $\sqrt{2}$. In fact, its elements are all of the form

$$x + y\sqrt{2} \quad \text{where } x \in \mathbf{Q} \text{ and } y \in \mathbf{Q}.$$

It's pretty obvious that most of the field axioms hold. The only one that's not obvious is the existence of reciprocals of nonzero elements, that is to say, the statement “ $(x + y\sqrt{2})^{-1}$ is of the form $x' + y'\sqrt{2}$ where x' and y' are rational and not both 0” is not so obvious. But the trick of “rationalizing the denominator” shows us how.

$$\frac{1}{x + y\sqrt{2}} = \frac{x - y\sqrt{2}}{(x + y\sqrt{2})(x - y\sqrt{2})} = \frac{x - y\sqrt{2}}{x^2 - 2y^2} = \frac{x}{x^2 - 2y^2} + \frac{-2y}{x^2 - 2y^2}\sqrt{2}$$

Note that $x^2 - 2y^2$ cannot be 0 when x and y are rational and not both 0. For if $x^2 - 2y^2 = 0$, then $2 = (x/y)^2$, and then $\sqrt{2}$ would be a rational number, which it isn't. Thus, $\mathbf{Q}(\sqrt{2})$ is a field.

The trick was to multiply and divide by the *conjugate*. Let's give a notation to this conjugate: $\overline{x + y\sqrt{2}} = x - y\sqrt{2}$. Conjugation has some nice properties. It preserves all the elements of the base field \mathbf{Q} , that is, if $x \in \mathbf{Q}$, then $\overline{x} = x$. It preserves addition and multiplication, that is, if α and β are elements of $\mathbf{Q}(\sqrt{2})$, then $\overline{\alpha + \beta} = \overline{\alpha} + \overline{\beta}$ and $\overline{\alpha\beta} = \overline{\alpha}\overline{\beta}$. Finally, the operation of conjugation, $\overline{} : \mathbf{Q}(\sqrt{2}) \rightarrow \mathbf{Q}(\sqrt{2})$, is its own inverse, $\overline{\overline{\alpha}} = \alpha$. Thus, conjugation is a field automorphism. Furthermore, the elements α it fixes, $\overline{\alpha} = \alpha$, are just the elements of the base field \mathbf{Q} .

2.3.2 The field of complex numbers \mathbf{C}

In the same way we just adjoined $\sqrt{2}$ to \mathbf{Q} to get $\mathbf{Q}(\sqrt{2})$, we can adjoin $\sqrt{-1}$ to \mathbf{R} to get $\mathbf{R}(\sqrt{-1})$, which is \mathbf{C} . Algebraically, the process is identical, but conceptually it's a little different because we thought that $\sqrt{2}$, being a real number, existed before we appended it to \mathbf{Q} , while it may not be so clear that $\sqrt{-1}$ exists before we append it to \mathbf{R} . But $\sqrt{-1}$, usually denoted i , has the property $i^2 = -1$, so it is an algebraic number since it's the root of the polynomial $x^2 + 1$. In fact, $\mathbf{R}(i)$ consists of elements of the form

$$x + yi \quad \text{with} \quad x, y \in \mathbf{R}$$

as described by Euler. Addition and subtraction are “coordinatewise”

$$(x_1 + y_1i) \pm (x_2 + y_2i) = (x_1 \pm x_2) + (y_1 \pm y_2)i$$

while multiplication is only slightly more complicated

$$\begin{aligned} (x_1 + y_1i)(x_2 + y_2i) &= x_1x_2 + x_1y_2i + x_2y_1i + y_1y_2i^2 \\ &= (x_1x_2 - y_1y_2) + (x_1y_2 + x_2y_1)i \end{aligned}$$

We can find reciprocals by rationalizing the denominator as we did above.

$$\frac{1}{x + yi} = \frac{x - yi}{(x + yi)(x - yi)} = \frac{x - yi}{x^2 + y^2} = \frac{x}{x^2 + y^2} + \frac{-y}{x^2 + y^2}i$$

We can define complex *conjugation* by $\overline{x + yi} = x - yi$. It's a field automorphism of \mathbf{C} , and its fixed subfield is \mathbf{R} .

We can also define a *norm* on \mathbf{C} once we have conjugation. For $z = x + yi \in \mathbf{C}$, let

$$|z|^2 = z\overline{z} = (x + yi)(x - yi) = x^2 + y^2.$$

Since $|z|^2$ is a nonnegative real number, it has a square root $|z|$.

A matrix representation of \mathbf{C} . Consider the subset C of the matrix ring $M_2(\mathbf{R})$ consisting of matrices of the form

$$\begin{bmatrix} x & y \\ -y & x \end{bmatrix} \quad \text{where} \quad x, y \in \mathbf{R}.$$

You can easily show that this is a subring of $M_2(\mathbf{R})$ since the 0 matrix and the identity matrix are of this form, the sum and difference of matrices of this form are of this form, and so is the product as you can see here

$$\begin{bmatrix} x & y \\ -y & x \end{bmatrix} \begin{bmatrix} u & v \\ -v & u \end{bmatrix} = \begin{bmatrix} xu - yv & xv + yu \\ -yu - vx & -yv + xu \end{bmatrix}.$$

Thus, C is a subring of $M_2(\mathbf{R})$. Furthermore, it's a commutative subring even though $M_2(\mathbf{R})$ is not a commutative ring since the same product results when the two factors are interchanged:

$$\begin{bmatrix} u & v \\ -v & u \end{bmatrix} \begin{bmatrix} x & y \\ -y & x \end{bmatrix} = \begin{bmatrix} ux - vy & uy + vx \\ -vx - uy & -vy + ux \end{bmatrix}.$$

Furthermore C is a field because nonzero matrices in it have inverses. For suppose not both x and y are 0. Then

$$\begin{bmatrix} x & y \\ -y & x \end{bmatrix} \begin{bmatrix} \frac{x}{x^2+y^2} & \frac{-y}{x^2+y^2} \\ \frac{y}{x^2+y^2} & \frac{x}{x^2+y^2} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

In fact, C is isomorphic to the complex field \mathbf{C} as described above. The isomorphism is described by the one-to-one correspondence

$$\begin{bmatrix} x & y \\ -y & x \end{bmatrix} \leftrightarrow x + yi.$$

Note that a real number x corresponds to the matrix $\begin{bmatrix} x & 0 \\ 0 & x \end{bmatrix}$ while a purely imaginary number yi corresponds to the matrix $\begin{bmatrix} 0 & y \\ -y & 0 \end{bmatrix}$.

Note that complex conjugation in this representation is just matrix transposition.

This alternate representation of the complex numbers as matrices directly explains how a complex number acts as a linear transformation on the real plane \mathbf{R}^2 . The complex number $x + yi$ maps a point (a, b) of \mathbf{R}^2 to the point $(ax + by, -ay + bx)$ since

$$\begin{bmatrix} x & y \\ -y & x \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} ax + by \\ -ay + bx \end{bmatrix}.$$

Matrix representations of various fields, rings, and groups are useful for two reasons. One is that they give us geometric interpretations for the elements as illustrated above. The other is that all the tools of linear algebra are available to us once we have the matrix representation.

2.3.3 General quadratic extensions

Now that we've seen a couple of quadratic extensions, let's see how it works in general.

Let F be a field and e an element of F that is not a square. In other words, the polynomial $x^2 - e$ has no roots in F . We'll consider ordered pairs $(a_1, a_2) \in F \times F$, but we'll write them as $a_1 + a_2\sqrt{e}$. We'll define addition coordinatewise

$$(a_1 + a_2\sqrt{e}) + (b_1 + b_2\sqrt{e}) = (a_1 + b_1) + (a_2 + b_2)\sqrt{e}$$

and define multiplication by

$$(a_1 + a_2\sqrt{e})(b_1 + b_2\sqrt{e}) = (a_1b_1 + ea_2b_2) + (a_1b_2 + a_2b_1)\sqrt{e}.$$

You can check that these definitions give us a ring. But, does it give us a field? As we did before, we'll find a reciprocal of a nonzero element $a_1 + a_2\sqrt{e}$

$$\frac{1}{a_1 + a_2\sqrt{e}} = \frac{a_1 - a_2\sqrt{e}}{(a_1 + a_2\sqrt{e})(a_1 - a_2\sqrt{e})} = \frac{a_1 - a_2\sqrt{e}}{a_1^2 - ea_2^2}$$

In order for this to be the reciprocal, all we have to do is show the denominator $a_1^2 - ea_2^2$ is not 0. In the case that $a_2 = 0$ we know $a_1 \neq 0$ since not both are 0, so in that case $a_1^2 - ea_2^2$ is not 0. That leaves us the case that $a_2 \neq 0$. Suppose that $a_1^2 - ea_2^2 = 0$. Then $ea_2^2 = a_1^2$, and dividing by a_2^2 , we conclude $e = (a_1/a_2)^2$. But e is not a square in F . Thus $a_1^2 - ea_2^2$ is not 0 in this case, too. Therefore, we've found the reciprocal.

Thus, we have a field, $F(\sqrt{e})$.

When we look at more general field extensions, we'll have a lot more theory, and we won't have details to check as we did here. That theory will involve the concept of "ideals" in a ring as discussed in section 3.6.

2.4 Real numbers and ordered fields

We'll look now at \mathbf{R} , the field of real numbers. What's so special about the real number field? For one thing, it's got an order on it; we can compare two real numbers x and y and say which is smaller or if they're equal. That's an extra structure on a field. We'll start by looking at this concept of ordered field.

Before we get too far, you should know that that isn't enough to distinguish \mathbf{R} from other fields. There are plenty of other ordered fields, such as \mathbf{Q} and all the fields between \mathbf{Q} and \mathbf{R} .

2.4.1 Ordered fields

The easiest way to define an ordered field is by saying it's partitioned into positive elements, negative elements, and 0, and requiring a couple properties on these parts.

Definition 2.24 (Ordered field). An *ordered field* consists of a field F along with a subset P whose elements are called *positive* such that

1. F is partitioned into three parts: P , $\{0\}$, and N where

$$N = \{x \in F \mid -x \in P\}$$

the elements of N are called *negative*;

2. the sum of two positive elements is positive; and
3. the product of two positive elements is positive.

Properties of ordered fields. You can show from this definition that

1. the sum of negative elements is negative
2. the product of a negative element and a positive element is negative
3. the product of two negative elements is positive
4. 1 is positive, and -1 is negative

Exercise 29. Prove the four properties above.

Examples. \mathbf{R} , \mathbf{Q} , and all fields between them are ordered fields where the usual positive numbers in the field form P .

Although \mathbf{Q} and \mathbf{R} are ordered fields, finite fields and \mathbf{C} have no ordering.

Exercise 30. Show that \mathbf{C} is not an ordered field. Hint: show why i can't be positive, zero, or negative.

The binary order relations. From P we can define the binary order relations $<$, \leq , $>$, and \geq . For instance, $x \leq y$ means $y - x$ is zero or positive, while $x < y$ means $y - x$ is positive. That can be stated formally as follows:

$$x \leq y \quad \text{iff} \quad y - x \in P \cup \{0\}$$

$$x < y \quad \text{iff} \quad y - x \in P.$$

All the expected properties of these order relations follow. Here are a few.

1. Trichotomy: For each pair x, y , exactly one of the three relations $x < y$, $x = y$, or $x > y$ holds.
2. Transitivity: $x < y$ and $y < z$ imply $x < z$.
3. If x is positive and $y < z$, then $xy < xz$.
4. If x is negative and $y < z$, then $xy > xz$.
5. If x is positive, then so is $1/x$.
6. For positive x and y , if $x < y$, then $1/y < 1/x$.

Exercise 31. Prove the six properties above.

Theorem 2.25. The characteristic of an ordered field is 0.

Proof. Suppose F is an ordered field of characteristic $p \neq 0$. Since 1 is positive, then any sum of 1s will be positive. Then p is positive. But p equals 0 which is not positive. A contradiction. Therefore an ordered field cannot have nonzero characteristic. Q.E.D.

It follows that \mathbf{Q} is a subfield of every ordered field.

Example 2.26. An ordered extension of the real numbers with infinite elements and infinitesimal elements.

We can give the field of rational functions $\mathbf{R}(x)$ an order as follows. First, we'll define when a polynomial $f(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$ is positive, and that will be when its leading coefficient a_n is a positive real number. Next, we'll define when a rational function $f(x)/g(x)$ is positive, and that will be when f and g are both positive polynomials or both negative polynomials. It follows that $f(x)/g(x)$ is negative one of f and g is positive and the other is negative. Only $0/g(x)$, which equals 0, won't be positive or negative. You can easily show that the sum and product of positive rational functions is positive.

The real numbers \mathbf{R} is an ordered subfield of $\mathbf{R}(x)$, meaning that it's a subfield and its elements have the same order whether the order on \mathbf{R} is used or the order on $\mathbf{R}(x)$ is used.

With this order, there are elements that are larger than any real number a , for example, $x > a$ since $x - a$ is positive. In other words, x is an infinite element. Likewise, there are positive elements that are smaller than any positive real number, $1/x$, for example, so $1/x$ is an infinitesimal number.

2.4.2 Archimedean orders

The last example is an example of an ordered field with infinite elements and infinitesimals. Every ordered field F is an extension of \mathbf{Q} , so we can define an infinite element of F to be an element $x \in F$ greater than every rational number, and we can define a positive infinitesimal element as a positive $x \in F$ smaller than every positive rational number. Note that the reciprocal of an infinite element is an infinitesimal, and vice versa.

Definition 2.27. An *Archimedean ordered field* or, more simply, an *Archimedean field*, is simply an ordered field F without infinite elements or infinitesimals.

Before Archimedes, Euclid used this property in his *Elements* in Book V and following books. The content of Book V is due to Eudoxus, so a better name for the Archimedean property would have been Eudoxus' property.

There are equivalent characteristics that could be used for the definition. Here are two. Each element of F is less than some integer. Each positive element of F is greater than the reciprocal of some positive integer.

Of course, the preceding example is a non-Archimedean field. Another interesting non-Archimedean field is that of surreal numbers created by John Conway. Surreal numbers include all real numbers, all ordinal numbers and more. Since ordinal numbers form a proper class, so do surreal numbers. For a nice introduction on surreal numbers, see Donald Knuth's book *Surreal Numbers*.

Still, there are loads of Archimedean fields, namely \mathbf{Q} , \mathbf{R} , and all the intermediate fields. We still haven't answered the question about what makes \mathbf{R} special. Before we go on, however, let's see how elements in an Archimedean field are determined by how they compare to rational numbers.

For an Archimedean field F , since F is ordered, it has characteristic 0, so it has as a subfield, indeed, an ordered subfield, the field of rational numbers \mathbf{Q} .

Theorem 2.28 (Density). Between any two distinct elements of an Archimedean field, there lies a rational number.

Proof. Let $x < y$ in an Archimedean field. We're looking for a rational number $\frac{m}{n}$ between x and y . If x is negative while y is positive, then the rational number 0 lies between them. We can reduce the case where they're both negative to the case where they're both positive by noting that if $\frac{m}{n}$ lies between $-x$ and $-y$, then $-\frac{m}{n}$ lies between x and y .

So we may assume that both x and y are positive. If we can find some multiple n of them so that $ny - nx > 1$, then some integer m lies between ny and nx , but $nx < m < ny$ gives $x < \frac{m}{n} < y$. And we can find such a multiple since $y - x$ is greater than the reciprocal $\frac{1}{n}$ of some positive integer since the field is Archimedean. Q.E.D.

An element a of F partitions \mathbf{Q} into two parts (L_a, R_a)

$$L_a = \{x \in \mathbf{Q} \mid x < a\} \quad \text{and} \quad R_a = \{x \in \mathbf{Q} \mid x \geq a\}.$$

These two parts have a special property.

Definition 2.29. A *Dedekind cut* of the rational numbers is a partition of \mathbf{Q} into two nonempty parts (L, R) —a left part L and a right part R —such that every element of L is less than every element of R . Furthermore, the left part does not have a greatest element.

Theorem 2.30. An element a of an Archimedean field F is determined by its Dedekind cut (L_a, R_a) . That is, if $(L_a, R_a) = (L_b, R_b)$, then $a = b$.

Proof. If $a \neq b$, then there is a rational number between them, so that rational number will be in one left part but the other right part. Q.E.D.

In an Archimedean field F not every Dedekind cut has to determine an element. For example, in \mathbf{Q} , the cut (L, R) where $L = \{x \mid x < 0 \text{ or } x^2 \leq 2\}$ and $R = \{x \mid x > 0 \text{ and } x^2 > 2\}$ is not the cut of any rational number. But that same cut with $\sqrt{2}$ included in \mathbf{R} is the cut of $\sqrt{2}$. The real numbers are special in that every cut is the cut of some real number.

Although there might not be a element of F for every cut, the cuts are enough to determine, along with the order on F and the field structure of \mathbf{Q} , the field structure of F .

It helps in proofs to cut in half the information of a Dedekind cut from (L, R) to just L . It is sufficient to define a Dedekind cut just in terms of of the left part. You can prove the following lemma to simplify the statement and the proof of the following theorem.

Lemma 2.31. If (L, R) is a Dedekind cut, then L has the following three properties

- i. L is a nonempty, proper subset of \mathbf{Q} ;
- ii. if $y \in L$ and $x \in \mathbf{Q}$ such that $x < y$, then $x \in L$; and
- iii. for each $x \in C$, there exists $y \in C$ such that $x < y$

Conversely, if L has these three properties, then (L, R) is a cut where R is the complement of L .

Theorem 2.32. In an Archimedean field F , addition and multiplication are determined by Dedekind cuts in the sense that If a and b are two elements of F , then the left part of their sum $a + b$ is determined by their left parts

$$L_{a+b} = \{x + y \mid x \in L_a \text{ and } y \in L_b\}.$$

If a and b are two positive elements of F , then the left part of their product is determined by their left parts

$$L_{ab} = \{xy \mid x \in L_a, x > 0, y \in L_b \text{ and } y > 0\} \cup \{x \mid x \leq 0\}.$$

2.4.3 Complete ordered fields

There are various definitions given for complete ordered fields, all logically equivalent. Here's one.

Definition 2.33. A *complete ordered field* is an Archimedean field that cannot be extended to a larger Archimedean field. Equivalently, every Dedekind cut determines an element of the field.

Completeness is the final property that characterizes \mathbf{R} . Actually, right now we haven't proved that there is *at least* one complete ordered field, and we haven't proved that there is *at most* one complete ordered field. Once we do, we can finally properly define \mathbf{R} .

Existence of a complete ordered field We'll start by stating the theorem which gives the components for one way of constructing a complete ordered field F . To make it complete, we just have to make sure that every Dedekind cut determines an element of the field. The way to do that, of course, to define the field to be the cuts, and the definition of the operations of addition and multiplication are determined by the cuts as seen in the last theorem.

Theorem 2.34. There is a complete ordered field F . Its elements are Dedekind cuts of \mathbf{Q} . If L_1 and L_2 are left parts of two cuts, then the left part of the sum is determined by the left part

$$L_+ = \{x + y \mid x \in L_1 \text{ and } y \in L_2\}.$$

If L is the left part a positive cut (one that contains at least one positive rational number), then its negation is determined by the left part

$$L_- = \{-x \mid x \notin L\}$$

except, if this L_- has a largest element, that largest element is removed. If L_1 and L_2 are left parts of two positive cuts, then the left part of the product is determined by the left part

$$L_\times = \{xy \mid x \in L_1, x > 0, y \in L_2 \text{ and } y > 0\} \cup \{x \mid x \leq 0\}.$$

There are many details to show to verify that R is a complete ordered field. First, that the sets L_+ , L_- , and L_\times are left parts. then the field axioms need to be verified, then the order axioms, then that's it's an Archimedean field. The last step, that it's complete is almost obvious from the construction. No one of these steps is difficult, but there are many details to check.

There are alternate ways to construct complete ordered fields. One is by means of Cauchy sequences. The spirit is different, but the result is the same, since, as we're about to see, there is only one complete ordered field.

Uniqueness of the complete ordered field We have to somehow exclude the possibility that there are two different Archimedean fields that can't be extended to larger Archimedean fields.

We don't want to count two isomorphic fields as being different, since, in essence, they're the same field but the names of the elements are just different. So, what we want is the following theorem.

Theorem 2.35. Any two complete ordered fields are isomorphic as ordered fields. Furthermore, there is only one isomorphism between them.

Proof. We may treat the field \mathbf{Q} as a subfield of the two complete ordered fields F_1 and F_2 . Then as a Dedekind cut determines an element $a_1 \in F_1$ and an element a_2 in F_2 , we have a bijection $F_1 \rightarrow F_2$. You only need to verify that preserves addition and multiplication, which it does, since in an Archimedean ring, addition and multiplication are determined by Dedekind cuts. Q.E.D.

R is the complete ordered field We now know that there is only one complete ordered field up to isomorphism. Any such complete ordered field may be taken as the real numbers.

2.5 Skew fields (division rings) and the quaternions

Sir William Rowan Hamilton, who early found that his road [to success with vectors] was obstructed—he knew not by what obstacle—so that many points which seemed within his reach were really inaccessible. He had done a considerable amount of good work, obstructed as he was, when, about the year 1843, he perceived clearly the obstruction to his progress in the shape of an old law which, prior to that time, had appeared like a law of common sense. The law in question is known as the *commutative* law of multiplication.

Kelland and Tait, 1873

2.5.1 Skew fields (division rings)

Skew fields, also called division rings, have all the properties of fields except that multiplication need not be commutative. When multiplication is not assumed to be commutative, a couple of the field axioms have to be stated in two forms, a left form and a right form. In particular, we require

1. there is a multiplicative identity, an element of F denoted 1, such that $\forall x, 1x = x = x1$;
2. there are multiplicative inverses of nonzero elements, that is, $\forall x \neq 0, \exists y, xy = 1 = yx$; and
3. multiplication distributes over addition, that is, $\forall x, \forall y, \forall z, x(y + z) = xy + xz$ and $\forall x, \forall y, \forall z, (y + z)x = yx + zx$.

All the other axioms remain the same, except we no longer require commutative multiplication.

The various properties of fields that follow from the field axioms also follow from the skew field axioms, although some have to be stated in two forms.

The most important skew field is the quaternions, mentioned next. Waring showed that there were no finite skew fields that weren't fields (a difficult proof).

2.5.2 The quaternions \mathbf{H}

We're not going to study skew fields, but one is of particular importance, the quaternions, denoted \mathbf{H} . The letter \mathbf{H} is in honor of Hamilton, their inventor.

We can define a quaternion a as an expression

$$a = a_0 + a_1i + a_2j + a_3k$$

where a_0, a_1, a_2 , and a_3 are real numbers and i, j , and k are formal symbols satisfying the properties

$$i^2 = j^2 = k^2 = -1$$

and

$$ij = k, jk = i, ki = j.$$

The i, j , and k are all square roots of -1 , but they don't commute as you can show from the definition that

$$ji = -k, kj = -i, ik = -j.$$

This doesn't lead to a commutative multiplication, but note that if a is real (i.e., its pure quaternion parts a_1, a_2 , and a_3 are all 0), then a will commute with any quaternion b .

Addition and subtraction are coordinatewise just like in \mathbf{C} . Here's multiplication.

$$\begin{aligned} & (a_0 + a_1i + a_2j + a_3k)(b_0 + b_1i + b_2j + b_3k) \\ &= (a_0b_0 - a_1b_1 - a_2b_2 - a_3b_3) \\ &+ (a_0b_1 + a_1b_0 + a_2b_3 - a_3b_2)i \\ &+ (a_0b_2 - a_1b_3 + a_2b_0 + a_3b_1)j \\ &+ (a_0b_3 + a_1b_2 - a_2b_1 - a_3b_0)k \end{aligned}$$

It's easy to check that all the axioms for a noncommutative ring are satisfied. The only thing left to in order to show that \mathbf{H} is a skew field is that reciprocals exist. We can use a variant of rationalizing the denominator to find the reciprocal of a quaternion.

$$\begin{aligned} \frac{1}{a_0 + a_1i + a_2j + a_3k} &= \frac{a_0 - a_1i - a_2j - a_3k}{(a_0 - a_1i - a_2j - a_3k)(a_0 + a_1i + a_2j + a_3k)} \\ &= \frac{a_0 - a_1i - a_2j - a_3k}{a_0^2 + a_1^2 + a_2^2 + a_3^2} \end{aligned}$$

Thus, a nonzero quaternion $a_0 + a_1i + a_2j + a_3k$, that is, one where not all of the real numbers a_0, a_1, a_2 , and a_3 are 0, has an inverse, since the denominator $a_0^2 + a_1^2 + a_2^2 + a_3^2$ is a nonzero real number.

The expression $a_0 - a_1i - a_2j - a_3k$ used to rationalize the denominator is the *conjugate* of the original quaternion $a_0 + a_1i + a_2j + a_3k$. It's worthwhile to have a notation for it.

$$\overline{a_0 + a_1i + a_2j + a_3k} = a_0 - a_1i - a_2j - a_3k,$$

as we do for \mathbf{C} . We'll also define the *norm* of a quaternion a by $|a|^2 = a\bar{a}$. It's a nonnegative real number, so it has a square root $|a|$. Note that $|a|^2 = a_0^2 + a_1^2 + a_2^2 + a_3^2$.

Thus, if a is a nonzero quaternion, then its inverse is $\frac{1}{a} = \frac{\bar{a}}{|a|^2}$.

For \mathbf{C} , the field of complex numbers, conjugation was a field automorphism, but for \mathbf{H} , it's not quite an automorphism. It has all of the properties of an automorphism except one. It preserves 0, 1, addition and subtraction $\overline{a \pm b} = \overline{a} \pm \overline{b}$, and reciprocation $\overline{1/a} = 1/\overline{a}$, but it reverses the order of multiplication $\overline{ab} = \overline{b} \overline{a}$. We'll call such a thing an *antiautomorphism*.

Note that \mathbf{H} extends \mathbf{C} in many ways. The assignment $x + iy \in \mathbf{C}$ to $x + iy \in \mathbf{H}$ is one, but $x + iy \in \mathbf{C}$ to $x + jy \in \mathbf{H}$ is another. There are, in fact, infinitely many ways that the skew field \mathbf{H} extends the field \mathbf{C} .

Theorem 2.36. The norm of a product is the product of the norms.

Proof. $|ab|^2 = ab\overline{ab} = ab\overline{b}\overline{a} = a|b|^2\overline{a} = a\overline{a}|b|^2 = |a|^2|b|^2$. Q.E.D.

If we unpack the equation $|a|^2|b|^2 = |ab|^2$, we'll get as a corollary Lagrange's identity on real numbers which shows how to express the product of two sums of four squares as the sum of four squares.

Corollary 2.37 (Lagrange). The product of the sum of four squares of integers is a sum of four squares of integers

$$\begin{aligned} & (a_0^2 + a_1^2 + a_2^2 + a_3^2)(b_0^2 + b_1^2 + b_2^2 + b_3^2) \\ &= (a_0b_0 - a_1b_1 - a_2b_2 - a_3b_3)^2 \\ &+ (a_0b_1 + a_1b_0 + a_2b_3 - a_3b_2)^2 \\ &+ (a_1b_2 + a_2b_1 + a_3b_1 - a_1b_3)^2 \\ &+ (a_2b_3 + a_3b_2 + a_1b_2 - a_2b_1)^2 \end{aligned}$$

Note that this equation not only works for real numbers, but also for integers, indeed when the coefficients lie in any commutative ring. Lagrange used this identity to show that every nonnegative integer n is the sum of four squares. The identity above is used to reduce the general case to the case when n is prime. Lagrange still had work to do to take care of the prime case.

Frobenius's theorem and the octonions. The quaternions are very special in the sense that they're the only finite-dimensional division algebra over \mathbf{R} other than \mathbf{R} itself and \mathbf{C} . This theorem was proved by Frobenius in 1877.

A division algebra over the real numbers R is a division ring (skew field) that has the reals as a subfield. Its dimension is the dimension it has as a vector space over \mathbf{R} .

There is also an eight-dimensional non-associative algebra over the real numbers called the *octonions*, \mathbf{O} . Octonions were discovered by John T. Graves in 1843. Although \mathbf{O} is not associative, it does satisfy weaker associativity identities when two of the three variables are the same: $x(xy) = (xx)y$, $x(yy) = (xy)y$, and $(xy)x = x(yx)$. It also satisfies the Moufang identities: $z(x(zy)) = ((zx)z)y$, $x(z(yz)) = ((xz)y)z$, $(zx)(yz) = (z(xy))z$, and $(zx)(yz) = z((xy)z)$. Furthermore \mathbf{O} has a norm.

Octonions over \mathbf{R} are a special case of a Cayley algebra over a field.

A matrix representation for \mathbf{H} . There are various matrix representations for \mathbf{H} . This one will make \mathbf{H} a subring of the real matrix ring $M_4(\mathbf{R})$. We'll represent 1 by the identity matrix, and i , j , and k by three other matrices which, you can verify, satisfy $i^2 = j^2 = k^2 = -1$ and $ij = k, jk = i, ki = j$.

$$1 \leftrightarrow \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad i \leftrightarrow \begin{bmatrix} 0 & -1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

$$j \leftrightarrow \begin{bmatrix} 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \end{bmatrix} \quad k \leftrightarrow \begin{bmatrix} 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

Then a generic quaternion $a + bi + cj + dk$ corresponds to the matrix

$$\begin{bmatrix} a & -b & -c & -d \\ b & a & -d & c \\ c & d & a & -b \\ d & -c & b & a \end{bmatrix}$$

Quaternions and geometry. Each quaternion a is the sum of a real part a_0 and a pure quaternion part $a_1i + a_2j + a_3k$. Hamilton called the real part a *scalar* and pure quaternion part a *vector*. We can interpret $a_1i + a_2j + a_3k$ as a vector $\mathbf{a} = (a_1, a_2, a_3)$ in \mathbf{R}^3 . Addition and subtraction of pure quaternions then are just ordinary vector addition and subtraction.

Hamilton recognized that the product of two vectors (pure quaternions) had both a vector component and a scalar component (the real part). The vector component of the product \mathbf{ab} of two pure quaternions Hamilton called the *vector product*, now often denoted $\mathbf{a} \times \mathbf{b}$ or $\mathbf{a} \vee \mathbf{b}$, and called the *cross product* or the *outer product*. The negation of the scalar component Hamilton called the *scalar product*, now often denoted $\mathbf{a} \cdot \mathbf{b}$, (\mathbf{a}, \mathbf{b}) , $\langle \mathbf{a}, \mathbf{b} \rangle$, or $\langle \mathbf{a} | \mathbf{b} \rangle$ and called the *dot product* or the *inner product*. Thus

$$\mathbf{ab} = \mathbf{a} \times \mathbf{b} - \mathbf{a} \cdot \mathbf{b}.$$

Hamilton's quaternions were very successful in the 19th century in the study of three-dimensional geometry.

Here's a typical problem from Kelland and Tait's 1873 *Introduction to Quaternions*. If three mutually perpendicular vectors be drawn from a point to a plane, the sum of the reciprocals of the squares of their lengths is independent of their directions.

Matrices were invented later in the 19th century. (But determinants were invented earlier!) Matrix algebra supplanted quaternion algebra in the early 20th century because (1) they described linear transformations, and (2) they weren't restricted to three dimensions.

Exercise 32. Show that \mathbf{H} can be represented as a subring of the complex matrix ring $M_2(\mathbf{C})$ where

$$1 \leftrightarrow \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad i \leftrightarrow \begin{bmatrix} i & 0 \\ 0 & -i \end{bmatrix}$$

$$j \leftrightarrow \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \quad k \leftrightarrow \begin{bmatrix} 0 & i \\ i & 0 \end{bmatrix}$$

so that a generic quaternion $a + bi + cj + dk$ corresponds to the matrix

$$\begin{bmatrix} a + bi & c + di \\ -c + di & a - bi \end{bmatrix}$$

Unit quaternions and S^3 as a group. The quaternions $a = a_0 + a_1i + a_2j + a_3k$ with norm 1 are called *unit quaternions*. Examples of unit quaternions are $\pm 1, \pm i, \pm j, \pm k$, but there are many more.

Unit quaternions are the quaternions for which $a_0^2 + a_1^2 + a_2^2 + a_3^2 = 1$. That equation is precisely the equation that defines the unit three sphere S^3 in 4-space \mathbf{R}^4 , although S^3 is usually described with different variables:

$$S^3 = \{(w, x, y, z) \in \mathbf{R}^4 \mid w^2 + x^2 + y^2 + z^2 = 1\}.$$

As we saw above, the product of the norms of two quaternions is the norm of the product, therefore multiplication is closed on this 3-sphere. Furthermore, 1 is a unit quaternion, and the reciprocal of a unit quaternion is another one, and, multiplication is associative, so multiplication of quaternions makes the 3-sphere S^3 into a group.

Chapter 3

Rings

Rings are things like \mathbf{Z} that have the three operations of addition, subtraction, and multiplication, but they don't need division. The lack of a division operation makes them more complicated and more interesting. The concept of *prime*, for example, is uninteresting for fields, but very interesting for \mathbf{Z} and other rings.

Most of our rings will have commutative multiplication, but some won't, so we won't require that multiplication be commutative in our definition. We will require that every ring have 1. The formal definition for rings is very similar to that for fields, but we leave out a couple of the requirements.

In this chapter we'll concentrate mainly on commutative rings and their properties. We'll consider commutative rings with various nice properties. Those rings with nice properties we'll give special names in increasing niceness such as *integral domain*, *unique factorization domain*, *principal ideal domain*, and *Euclidean domain*.

3.1 Introduction to rings

A *ring* is a set equipped with two binary operations, one called *addition* and the other called *multiplication*, denoted in the usual manner, which are both associative, addition is commutative, both have identity elements (the additive identity denoted 0 and the multiplicative identity denoted 1), addition has inverse elements (the inverse of x denoted $-x$), and multiplication distributes over addition. If, furthermore, multiplication is commutative, then the ring is called a *commutative ring*.

3.1.1 Definition and properties of rings

Here is a more complete definition.

Definition 3.1. A *ring* R consists of

1. a set, also denoted R and called the *underlying set* of the ring;
2. a binary operation $+$: $R \times R \rightarrow R$ called *addition*, which maps an ordered pair $(x, y) \in R \times R$ to its *sum* denoted $x + y$;

3. another binary operation $\cdot : R \times R \rightarrow R$ called *multiplication*, which maps an ordered pair $(x, y) \in R \times R$ to its *product* denoted $x \cdot y$, or more simply just xy ; such that
4. addition is commutative, that is, $\forall x, \forall y, x + y = y + x$;
5. addition is associative, that is, $\forall x, \forall y, (x + y) + z = x + (y + z)$;
6. multiplication is associative, that is, $\forall x, \forall y, (xy)z = x(yz)$;
7. there is an additive identity, an element of F denoted 0 , such that $\forall x, 0 + x = x$;
8. there is a multiplicative identity, an element of F denoted 1 , such that $\forall x, 1x = x$;
9. there are additive inverses, that is, $\forall x, \exists y, x + y = 0$; and
10. multiplication distributes over addition, that is, $\forall x, \forall y, \forall z, x(y + z) = xy + xz$.

When multiplication is also commutative, that is, $\forall x, \forall y, xy = yx$, the ring is called a *commutative ring*. The conditions for a ring are often call the *ring axioms*.

Subtraction, multiples, and powers. As we did with fields, we can define subtraction, integral multiples, and nonnegative integral powers. We won't have division or negative integral powers since we don't have reciprocals.

As before, we define subtraction in terms of negation. The *difference* of two elements x and y is $x - y = x + (-y)$. The expected properties of subtraction all follow from the ring axioms. For instance, multiplication distributes over subtraction.

Likewise, we can define integral multiples of elements in a ring. Define $0x$ as 0 , then inductively define $(n + 1)x = x + nx$ when $n \geq 0$. Then if $-n$ is a negative integer, define $-nx$ as $-(nx)$. The usual properties of multiples, like $(m + n)x = mx + nx$ still hold.

Furthermore, we can define positive integral powers of x . Define x^1 as x for a base case, and inductively, $x^{n+1} = xx^n$. Thus nx is the product of n x 's. For instance, $x^3 = xxx$. Since rings needn't have reciprocals, we can't define negative integral powers of x .

Examples 3.2 (rings). Of course, all fields are automatically rings, but what are some other rings? We've talked about some others already, including

1. the ring of integers \mathbf{Z} which includes all integers (whole numbers)—positive, negative, or 0 .
2. the ring of polynomials $R[x]$ with coefficients in a commutative ring R .
3. the matrix ring $M_n(R)$ of $n \times n$ matrices with entries in a commutative ring R . This example is a noncommutative ring when $n \geq 2$.
4. the ring of upper triangular matrices is a subring of $M_n(R)$.
5. the cyclic ring \mathbf{Z}_n , the ring of integers modulo n , where n is a particular integer.
6. the powerset $\mathcal{P}(S)$ consisting of subsets of a set S becomes a ring, called a Boolean ring, where $A + B$ is the symmetric difference and AB is the intersection of two subsets A and B .

Properties that follow from the ring axioms. There are numerous useful properties that follow from the axioms, but not so many as follow from the field axioms. Here's a list of several of them.

1. 0 is unique. That is, there is only one element x of a ring that has the property that $\forall y, x + y = y$. Likewise, 1 is unique.
2. Multiplication distributes over subtraction. $x(y - z) = xy - xz$ and $(y - z)x = yx - zx$.
3. $-0 = 0$.
4. $0x = 0$.
5. $(-1)x = -x$, $(-x)y = -(xy) = x(-y)$, and $(-x)(-y) = xy$.

There are some expected properties that are not included here. I'll show why not using examples from \mathbf{Z}_6 .

1. If the product of two elements is 0, $xy = 0$, it does not follow that either $x = 0$ or $y = 0$. For example, in \mathbf{Z}_6 the product of 2 and 3 is 0.
2. Cancellation does not always work. That is, if $xy = xz$ and $x \neq 0$, it doesn't follow that $y = z$. For example, in \mathbf{Z}_6 , $3 \cdot 2 = 3 \cdot 4$, but $2 \neq 4$.

3.1.2 Products of rings

If R_1 and R_2 are two rings, you can construct their product ring R . The underlying set of R is the product $R_1 \times R_2$ of the underlying sets of the two rings, and addition, subtraction, and multiplication are coordinatewise. Thus,

$$(x_1, x_2) \pm (y_1, y_2) = (x_1 \pm y_1, x_2 \pm y_2) \quad \text{and} \quad (x_1, x_2)(y_1, y_2) = (x_1y_1, x_2y_2).$$

The additive identity in $R_1 \times R_2$ is $0 = (0, 0)$, and the multiplicative identity is $1 = (1, 1)$. Since all the operations are performed coordinatewise, the ring axioms are satisfied in $R_1 \times R_2$, so it's a ring.

The projection functions $\pi_1 : R_1 \times R_2 \rightarrow R_1$ and $\pi_2 : R_1 \times R_2 \rightarrow R_2$ defined by $\pi_1(x_1, x_2) = x_1$ and $\pi_2(x_1, x_2) = x_2$ are both ring homomorphisms. They preserve addition, multiplication, and 1.

Products of more than 2 rings can be defined analogously, even products of infinitely many rings.

We didn't discuss products of fields in the chapter on field because the product of two fields is not another field. It is at least a ring, however.

3.1.3 Integral domains

Much of the time we will want the cancellation property that was mentioned above to hold, so we'll give a special name to commutative rings that satisfy them. It will help if we make a couple of definitions.

Definition 3.3. A nonzero element x in a commutative ring is a *zero-divisor* if there exists a nonzero y such that $xy = 0$. Of course, 0 is always a zero-divisor. We'll say a commutative ring *has no zero divisors* if 0 is the only zero-divisor.

Definition 3.4. We'll say a commutative ring satisfies the *cancellation law* if

$$\forall x \neq 0, \forall y, \forall z, xy = xz \text{ implies } y = z.$$

We found in the example above that 2 and 3 are zero-divisors in \mathbf{Z}_6 , and that \mathbf{Z}_6 did not satisfy the cancellation law. You can examine \mathbf{Z}_n to determine which nonzero elements are zero-divisors and which have reciprocals.

There's a connection between zero-divisors and the cancellation law.

Theorem 3.5. A commutative ring satisfies the cancellation law if and only if it has no zero-divisors.

Proof. Suppose the ring satisfies the cancellation law. Let x be a nonzero element in the ring. If $xy = 0$, then $xy = x0$, so by that cancellation law, $y = 0$. Then x can't be a zero-divisor. Thus the ring has no zero-divisors.

Next suppose that the ring has no zero-divisors. We'll show it satisfies the cancellation law. If $x \neq 0$ and $xy = xz$, then $x(y - z) = 0$, and since x is not a zero divisor, therefore $y - z = 0$, so $y = z$. Thus the ring satisfies the cancellation law. Q.E.D.

Definition 3.6 (integral domain). An *integral domain* is a commutative ring D in which $0 \neq 1$ that satisfies one of the two equivalent conditions: it has no zero-divisors, or it satisfies the cancellation law.

All the fields and most of the examples of commutative rings we've looked at are integral domains, but \mathbf{Z}_n is not an integral domain if n is not a prime number.

Note that any subring of a field or an integral domain will be an integral domain since the subring still won't have any zero-divisors.

Note that products of (nontrivial) rings are never integral domains since they always have the zero divisors $(1, 0)$ and $(0, 1)$ whose product is 0.

Corollary 2.13 stated that the characteristic of a field is either 0 or a prime number. The proof there works as well for integral domains. The characteristic of an integral domain is either 0 or a prime number.

Group rings You can form a ring $\mathbf{Z}G$ out of a group G as follows. Assume that G is written multiplicatively. The finite formal sums of elements of G are the elements of $\mathbf{Z}G$. Thus, if n is a nonnegative integer and $a_1, \dots, a_n \in G$, then the formal sum $x_1a_1 + \dots + x_na_n$ names an element of the group ring $\mathbf{Z}G$. Addition is coordinatewise. Multiplication uses the group operation.

This definition can be generalized so that group rings have their coordinates in any commutative ring R , not just \mathbf{Z} . This results in a group ring RG .

Exercise 33. Let G be the two element cyclic group $G = \{1, a\}$ where $a^2 = 1$. A typical element of $\mathbf{Z}G$ is $x + ya$ where $x, y \in \mathbf{Z}$. Multiplication is defined by $(x_1 + y_1a)(x_2 + y_2a) = (x_1x_2 + y_1y_2) + (x_1y_2 + x_2y_1)a$. Show that the square of any nonzero element in $\mathbf{Z}G$ is not zero, but show that $\mathbf{Z}G$ does have zero-divisors by finding a pair of nonzero elements whose product is 0.

3.1.4 The Gaussian integers, $\mathbf{Z}[i]$

One important example of an integral domain is that of the Gaussian integers $\mathbf{Z}[i]$. Its elements are of the form $x + yi$ where $x, y \in \mathbf{Z}$, so they can be viewed as a lattice of points in the complex plane as in figure 3.1. You can check that $\mathbf{Z}[i]$ is closed under addition, subtraction, multiplication, and includes 1, so it is a subring of the field \mathbf{C} . Therefore, it's an integral domain. We'll see later that $\mathbf{Z}[i]$ is a particularly nice integral domain called a Euclidean domain.

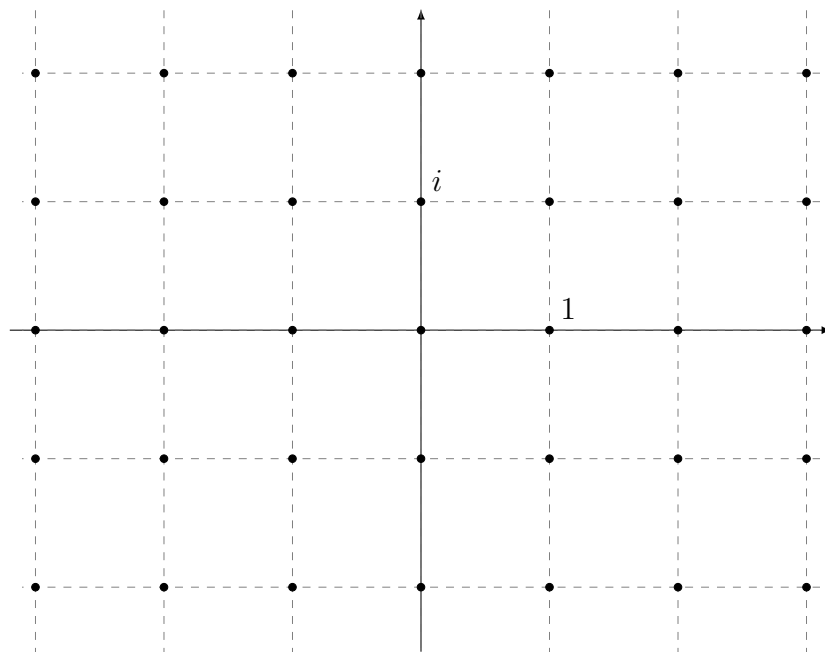


Figure 3.1: Lattice of Gaussian integers $\mathbf{Z}[i]$

There are four units (elements having reciprocals) in the Gaussian integers. Besides 1 and -1 , i and $-i$ are also units. Note that $(1 + i)(1 - i) = 2$, so 2 is not prime in $\mathbf{Z}[i]$ even though it is prime in \mathbf{Z} .

We'll come back to $\mathbf{Z}[i]$ when we study Euclidean domains in section 3.8.4. Also $\mathbf{Z}[i]$ is an example of a “ring of integers” to be defined in section 3.11.

Eisenstein integers. The Eisenstein integers are similar to the Gaussian integers, but instead of consisting of a square lattices of complex numbers, they consist of a triangular lattice of complex numbers. They include complex numbers of the form $z = x + y\omega$ where ω is the cube root of 1, $\omega = \frac{1}{2}(-1 + i\sqrt{3}) = e^{2\pi i/3}$. See figure 3.3 for the lattice of Eisenstein integers.

3.1.5 Finite fields again

We won't find any examples of finite integral domains that aren't fields because there aren't any.

Theorem 3.7 (Wedderburn). If R is a finite integral domain, then R is a field.

Proof. Let x be a nonzero element of R . Consider the positive powers of x :

$$x, x^2, x^3, \dots, x^n \dots$$

Since there are infinitely many powers, but only finitely many elements in R , therefore at least two distinct powers are equal. Let, then, $x^m = x^n$ with $m < n$. Cancel x^m from each side of the equation (which is possible because R is an integral domain) to conclude $x^{n-m} = 1$. Therefore, the reciprocal of x is x^{n-m-1} . Therefore, every nonzero element has an inverse. Q.E.D.

This theorem can be used to give a short proof that \mathbf{Z}_p is a field when p is a prime, since it's easy to show that \mathbf{Z}_p is an integral domain. We'll show it has no zero-divisors. Suppose that $xy \equiv 0 \pmod{p}$. Then $p|xy$. But if a prime divides a product, it divides one of the factors, so either $p|x$ or $p|y$, in other words, either $x \equiv 0 \pmod{p}$ or $y \equiv 0 \pmod{p}$. Thus, \mathbf{Z}_p is an integral domain, and hence, by the above theorem, it's a field.

Our earlier, more complicated proof used the extended Euclidean algorithm to find an inverse for x . That's actually a much more efficient way to find the inverse than to look through the powers of x .

3.2 Factoring \mathbf{Z}_n by the Chinese remainder theorem

We'll look at the structure of the cyclic ring \mathbf{Z}_n when n is composite in more detail. In particular, when n is not a power of a prime number, then \mathbf{Z}_n is a product of smaller cyclic rings.

3.2.1 The Chinese remainder theorem

This theorem says that if m and k are relatively prime and $n = mk$, then $\mathbf{Z}_n \cong \mathbf{Z}_m \times \mathbf{Z}_k$. Let's illustrate that with $m = 7$ and $k = 12$ to show how $\mathbf{Z}_{84} \cong \mathbf{Z}_7 \times \mathbf{Z}_{12}$. Starting with a number x modulo 84, we'll get a pair of numbers, one being x modulo 7, the other x modulo 12. We can display this in a 7×12 table where each row is a number modulo 7, each column a number modulo 12, and the entry at row i and column j is that number which is i modulo 7 and j modulo 12.

It's easy to construct the table. Start filling the diagonal. After you reach the last row, go next to the top row, and after you reach the right column, go next to the left column.

	0	1	2	3	4	5	6	7	8	9	10	11
0	0	49	14	63	28	77	42	7	56	21	70	35
1	36	1	50	15	64	29	78	43	8	57	22	71
2	72	37	2	51	16	65	30	79	44	9	58	23
3	24	73	38	3	52	17	66	31	80	45	10	59
4	60	25	74	39	4	53	18	67	32	81	46	11
5	12	61	26	75	40	5	54	19	68	33	82	47
6	48	13	62	27	76	41	6	55	20	69	34	83

All the numbers in the first row are congruent to 0 modulo 7, so they're divisible by 7, but looking at them, they seem to be rather randomly arranged. Likewise, all the numbers in the first column are divisible by 12.

The pair of linear congruences $x \equiv i \pmod{7}$ and $x \equiv j \pmod{12}$ can be easily solved for x in by looking in row i and column j .

For example, take this Chinese remainder problem. Find a number such that when you divide it by 7 you get a remainder of 3, but when you divide it by 12 you get a remainder of 8. The answer, 80, is right in the table.

Theorem 3.8 (Chinese remainder theorem). Suppose that $n = km$ where k and m are relatively prime. Then

$$\mathbf{Z}_n \cong \mathbf{Z}_k \times \mathbf{Z}_m.$$

More generally, if n is the product $k_1 \cdots k_r$ where the factors are pairwise relatively prime, then

$$\mathbf{Z}_n \cong \mathbf{Z}_{k_1} \times \cdots \times \mathbf{Z}_{k_r} = \prod_{i=1}^r \mathbf{Z}_{k_i}.$$

In particular, if the prime factorization of n is $n = p_1^{e_1} \cdots p_r^{e_r}$. Then the cyclic ring \mathbf{Z}_n factors as the product of the cyclic rings $\mathbf{Z}_{p_i^{e_i}}$, that is,

$$\mathbf{Z}_n \cong \prod_{i=1}^r \mathbf{Z}_{p_i^{e_i}}.$$

Proof. The third statement is a special case of the second.

The second follows from the first by induction on r .

That leaves us with the first statement. In one direction, $\mathbf{Z}_n \rightarrow \mathbf{Z}_k \times \mathbf{Z}_m$, the function giving the isomorphism is fairly obvious; it's built of the two functions $\mathbf{Z}_n \rightarrow \mathbf{Z}_k$ and $\mathbf{Z}_n \rightarrow \mathbf{Z}_m$ that are easy to describe.

There is an obvious candidate for a ring function $\mathbf{Z}_n \rightarrow \mathbf{Z}_k$, namely $[x]_n \mapsto [x]_k$ by which is meant the equivalence class of x modulo n is sent to the equivalence class of x modulo k .

First, we have to check that this function is well defined. Suppose $[x]_n = [y]_n$. Then $x \equiv y \pmod{n}$, so $n|(x - y)$. But $k|n$, therefore $k|(x - y)$. Hence, $x \equiv y \pmod{k}$, and $[x]_k = [y]_k$. So the function is well-defined.

You can check the rest, that this function preserves the ring operation so that it's a ring homomorphism.

Putting together the two ring homomorphisms $\mathbf{Z}_n \rightarrow \mathbf{Z}_k$ and $\mathbf{Z}_n \rightarrow \mathbf{Z}_m$ we have a ring homomorphism

$$\begin{aligned} \mathbf{Z}_n &\rightarrow \mathbf{Z}_k \times \mathbf{Z}_m \\ [x]_n &\mapsto ([x]_k, [x]_m) \end{aligned}$$

In order to show that this is an isomorphism, all we need to do is to show that it's a bijection, and for that, all we need to do is to show that it's an injection since the sets \mathbf{Z}_n and $\mathbf{Z}_k \times \mathbf{Z}_m$ have the same cardinality.

Suppose that $[x]_n$ and $[y]_n$ are sent to the same element in $\mathbf{Z}_k \times \mathbf{Z}_m$. Then $[x]_k = [y]_k$ and $[x]_m = [y]_m$, that is, $k|(x - y)$ and $m|(x - y)$. Since they both divide $x - y$, so does their least common multiple. But they're relatively prime, so their LCM is their product, n . Thus $n|(x - y)$, so $[x]_n = [y]_n$. Therefore, this is a one-to-one function, hence a one-to-one correspondence. Thus, the ring homomorphism is an isomorphism. Q.E.D.

The inverse. Well, since it's a bijection, it shouldn't be too hard to find its inverse $\mathbf{Z}_k \times \mathbf{Z}_m \rightarrow \mathbf{Z}_n$. In other words, solve for $x \pmod{n}$ the pair of simultaneous congruences

$$\begin{aligned}x &\equiv a \pmod{k} \\x &\equiv b \pmod{m}\end{aligned}$$

It's too much work to construct the entire $k \times m$ table as was done for the 7×12 . There's a better way.

We can find a solution with the extended Euclidean algorithm. Since $\text{gcd}(m, k) = 1$, therefore 1 is a linear combination of m and k , that is, there are integers s and t so that $sm + tk = 1$. Multiply by $b - a$ to conclude $s(b - a)m + t(b - a)k = b - a$. Therefore, $t(b - a)k + a = b - s(b - a)m$. Let that be x . Then $x \equiv a \pmod{k}$ and $x \equiv b \pmod{m}$ as required.

Problems like this in indeterminate analysis were solved in ancient China and in ancient India. The earliest appeared in *Sunzi suanjing* (*Master Sun's Mathematical Manual*) in the about the fourth century C.E. in China. In 1247 Qin Jiushao gave a general method for solving linear congruences in his *Shushu jiuzhang* (*Mathematical Treatise in Nine Sections*).

3.2.2 Brahmagupta's solution

In India in the seventh century C.E., Brahmagupta also gave a general algorithm for solving these linear congruences in his *Brāhmasphuṭasiddhānta* (*Correct Astronomical System of Brahma*). If more than two congruences were given, he first reduced the problem to solving pairs of congruences as we did above. His solution is the one described above.

As an example, find $x \pmod{210}$ if

$$\begin{aligned}x &\equiv 11 \pmod{45} \\x &\equiv 4 \pmod{56}\end{aligned}$$

Here's how he did it in modern notation, explained with the numerical example above.

We're looking for a value of x so that $x = 45s + 11 = 56t + 4$ for some integers s and t . So we need s and t so that $45s + 7 = 56t$. That reduces to $45(s - t) + 7 = 11t$. Let $s' = s - t$. To solve $45s' + 7 = 11t$, since $45 = 4 \cdot 11 + 1$, reduce it to $s' + 7 = 11(t - 4s')$. Let $t' = t - 4s'$. We can solve $s' + 7 = 11t'$ by setting $s' = 4$ and $t' = 1$. Substituting these in the defining equations, we find $t = t' + 4s' = 17$, and $s = s' + t = 21$. Therefore, $x = 45s + 11 = 956$, the answer.

Of course, Brahmagupta did not use variables. His is solution was described as a fairly simple algorithm that just used the four arithmetic operations.

3.2.3 Qin Jiushao's solution

The algorithm that Qin Jiushao described was fairly different and applied directly to many linear congruences so long as the moduli were pairwise relatively prime. Let's illustrate it with the system of three congruences

$$\begin{aligned}x &\equiv 45 \pmod{121} \\x &\equiv 31 \pmod{63} \\x &\equiv 30 \pmod{100}\end{aligned}$$

Since the moduli are pairwise relatively prime, we can find a unique solution to this system modulo 762300, the product of the moduli.

Step 1. For each modulus, find a reciprocal of the product of the remaining moduli modulo the given modulus. For the first modulus, 121, that means we need the reciprocal of 6300 modulo 121, that is, we need to solve

$$6300y \equiv 1 \pmod{121}.$$

That's the same as $8y \equiv 1 \pmod{121}$. The extended Euclidean algorithm gives us $1 = (-15) \cdot 8 + 1 \cdot 121$, so $y = -15$ is a solution.

For the second modulus, 63, we need the reciprocal of 12100 modulo 63. That's the same as the reciprocal of 4 modulo 63, which is 16.

For the third modulus, 100, we need the reciprocal of 7623 modulo 100. That's the same as the reciprocal of 23 modulo 100. The Chinese mathematicians called finding a reciprocal modulo n by the term "finding one". By the extended Euclidean algorithm, $(-13) \cdot 23 + 3 \cdot 8 = 1$, so -13 is the reciprocal of 23 modulo 100.

Step 2. To get x sum three products abc , one for each congruence, where a is the constant in the congruence, b is the product of the other moduli, and c is the reciprocal found in the previous step. That gives us

$$\begin{aligned} & 45 \cdot 6300 \cdot (-15) \\ & + 31 \cdot 12100 \cdot 16 \\ & + 30 \cdot 7623 \cdot (-13) \\ & = -283515 + 6001600 - 2972970 = 2745115 \end{aligned}$$

and then reduce this number modulo the product 762300 of all three moduli. That gives a final answer of $x \equiv 458215 \pmod{762300}$.

Exercise 34. Solve the following system of simultaneous linear congruences. You can use either Brahmagupta's algorithm, Qin Jiushao's algorithm, or something of your own devising.

$$\begin{aligned} x &\equiv 4 \pmod{33} \\ x &\equiv 22 \pmod{35} \\ x &\equiv 41 \pmod{53} \end{aligned}$$

Be sure to show how you derived the solution.

3.3 Boolean rings

Representing by x the class "men," and by y "Asiatics," let z represent the adjective "white" to the collection of men expressed by the phrase "Men except Asiatics," is the same as to say "White men except white Asiatics." Hence we have

$$z(x - y) + zx - zy.$$

This is also in accordance with the laws of ordinary algebra.

George Boole, 1854. *An Investigation of the Laws of Thought on which are founded the mathematical theories of logic and probabilities.*

George Boole (1815-1864). Boole wanted to bring logic into the realm of mathematics, which he did by algebrizing it.

We'll incorporate his investigations in our study of ring theory, but change his notation slightly. Boole did not allow a sum of two things unless they were disjoint, so $x + x$ had no meaning for him. We'll just take $+$ to be an exclusive or (symmetric difference), so $x + x$ will be 0 for us.

3.3.1 Introduction to Boolean rings

We saw before that powerset $\mathcal{P}(S)$ of a set S becomes a ring when we define $A + B$ to be the symmetric difference and AB to be the intersection of two subsets A and B . The 0 element of the ring is the emptyset \emptyset , while the 1 element is S . The complement of a subset A is $1 - A$ (which equals $1 + A$).

We'll define what a Boolean ring is in terms of idempotents.

Definition 3.9. An element e of a ring is said to be *idempotent* when $e^2 = e$.

Notice that 0 and 1 are always idempotent in any ring.

Other examples of idempotent elements in rings are projections. The transformation $f : \mathbf{R}^2 \rightarrow \mathbf{R}^2$ which projects a point in the plane to the x -axis, defined by $f(x, y) = (x, 0)$, is idempotent as is any projection from a space to a subspace of itself.

Definition 3.10. A *Boolean ring* is a ring in which every element is idempotent.

The ring $\mathcal{P}(S)$ is evidently an example of a Boolean ring.

Two properties that follow from the definition are (1) that the a Boolean ring has characteristic 2, (2) Boolean rings are commutative.

Theorem 3.11. A nontrivial Boolean ring has characteristic 2.

Proof. Since $1 + 1$ is idempotent, $(1 + 1)^2 = 1 + 1$. Therefore, $1 + 1 + 1 + 1 = 1 + 1$, and so $1 + 1 = 0$. Q.E.D.

As in any ring of characteristic 2, negation does nothing, $-x = x$, and subtraction is the same as addition, $x - y = x + y$.

Theorem 3.12. Boolean rings are commutative.

Proof. Let x and y be two elements of a Boolean ring. Since $x + y$ is idempotent, $(x + y)^2 = x + y$. Expanding that equation using the fact that multiplication distributes over addition in every ring, commutative or not, it follows that $x^2 + xy + yx + y^2 = x + y$. But $x^2 = x$ and $y^2 = y$, so that last equation simplifies to $xy + yx = 0$. Therefore, $xy = -yx$, and $-yx = yx$, so $xy = yx$. Q.E.D.

Boolean rings are the same thing as something called Boolean algebras, but the approaches are different. A Boolean ring is thought of as a special kind of ring, while a Boolean algebra is a special kind of partially ordered set whose elements are truth values. Boolean algebras are reviewed in the appendix section [A.3.3](#).

Table [3.1](#) compares common notations in Boolean algebras, set theory, and Boolean rings. Here, P and Q are propositions or predicates, A and B are subsets of a set Ω , and x and y

Boolean algebras	Set theory	Boolean rings
T (true)	Ω	1
F (false)	\emptyset	0
$P \wedge Q$ (and)	$A \cap B$	xy
$P \vee Q$ (inclusive or)	$A \cup B$	$x + y + xy$
$P \oplus Q$ (exclusive or)	$A \oplus B$	$x + y$
$\neg P$ (not)	A^c	$1 + x$
$P \iff Q$	$A = B$	$x = y$
$P \implies Q$	$A \subseteq B$	$xy = x$
$T \vee Q \iff T$	$\Omega \cup B = \Omega$	$1 + y + 1y = 1$
$F \vee Q \iff Q$	$\emptyset \cup B = B$	$0 + y + 0y = y$
$T \wedge Q \iff Q$	$\Omega \cap B = B$	$1y = y$
$F \wedge Q \iff F$	$\emptyset \cap B = \emptyset$	$0y = 0$
$P \wedge Q \iff Q \wedge P$	$P \cap Q = Q \cap P$	$xy = yx$
$P \vee Q \iff Q \vee P$	$P \cup Q = Q \cup P$	$x + y + xy = y + x + yx$
$\neg(P \wedge Q) \iff \neg P \vee \neg Q$	$(A \cap B)^c = A^c \cup B^c$	$1 + xy = (1 + x) + (1 + y) + (1 + x)(1 + y)$
$\neg(P \vee Q) \iff \neg P \wedge \neg Q$	$(A \cup B)^c = A^c \cap B^c$	$1 + (x + y + xy) = (1 + x)(1 + y)$

Table 3.1: Notations in Boolean algebras, set theory, and Boolean rings.

are elements of a Boolean ring. These are just a few correspondences. You can add many more.

For more on Boolean algebras, see section A.3.3 in the appendix.

Free Boolean rings Some Boolean rings, called free Boolean rings, have special properties. Given a set S whose elements are called *generators*, the *free Boolean ring* on S is the Boolean ring $B(S)$ which comes equipped with a function $\iota : S \rightarrow B(S)$ that satisfies the following universal property: for each Boolean ring R and function $f : S \rightarrow R$, there exists a unique ring homomorphism $\hat{f} : B(S) \rightarrow R$ such that $f \circ \iota = \hat{f}$.

Examples 3.13 (Free Boolean rings). If $S = \emptyset$ is the empty set, then $B(\emptyset)$ consists of only two elements, 0 and 1, which when identified with truth values are \perp and \top , respectively.

If $S = \{p\}$ is a singleton set, then $B(\{p\})$ has four elements, 0, p , $1 - p$, and 1, which when identified with truth values are \perp , p , $\neg p$, and \top .

If $S = \{p, q\}$ has two elements, then $B(\{p, q\})$ has 16 elements. They are displayed in figure 3.2 as truth values. It is isomorphic to the Boolean ring which is the powerset of a set of four elements in figure A.1.

3.3.2 Factoring Boolean rings

Suppose that a set S is partitioned into subsets S_1, S_2, \dots, S_n . That means S is the union of all these subsets, and they are pairwise disjoint. Then the ring $\wp(S)$ is isomorphic to a product of the rings $\wp(S_i)$. The function

$$\begin{aligned} \wp(S) &\cong \wp(S_1) \times \wp(S_2) \times \dots \times \wp(S_n) \\ A &\mapsto (A \cap S_1, A \cap S_2, \dots, A \cap S_n) \end{aligned}$$

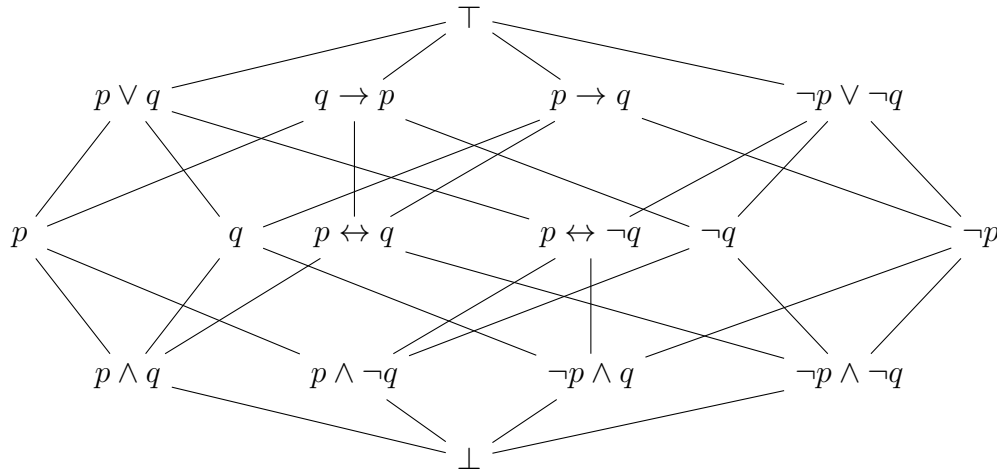


Figure 3.2: Free Boolean ring on two elements

gives the ring homomorphism in one direction, and it's a bijection since A is the disjoint union of the terms on the right.

In fact, this works even when S is partitioned into arbitrarily many subsets. Since S is the disjoint union of its singletons $S = \cup_{x \in S} \{x\}$, therefore $\mathcal{P} = \prod_{x \in S} \mathcal{P}(\{x\})$. In other words, \mathcal{P} is a power of the 2-element ring.

Factoring works in a general Boolean ring as well as those of the form $\mathcal{P}(S)$. Let R be a Boolean ring, and e any idempotent in it other than 0 or 1. Let $\bar{e} = 1 - e$, so that $1 = e + \bar{e}$ from which it follows that $x = xe + x\bar{e}$ for all $x \in R$. Let $R_e = \{xe \mid x \in R\}$, and let $R_{\bar{e}} = \{x\bar{e} \mid x \in R\}$. You can check that both R_e and $R_{\bar{e}}$ are Boolean rings, where the multiplicative identities are e and \bar{e} , respectively. Furthermore,

$$\begin{aligned} R &\cong R_e \times R_{\bar{e}} \\ x &\mapsto (xe, x\bar{e}) \end{aligned}$$

3.3.3 A partial order on a Boolean ring

If we define $x \preceq y$ to mean $xy = x$, then our Boolean ring will have a partial ordering.

Recall that a partial ordering \preceq on a set is a reflexive, antisymmetric, and transitive relation.

1. Reflexive: $x \preceq x$, since $x^2 = x$.
2. Antisymmetric: $x \preceq y$ and $y \preceq x$ imply $x = y$, since $xy = x$ and $yx = y$ imply $x = y$.
3. Transitive: $x \preceq y$ and $y \preceq z$ imply $x \preceq z$, since $xy = x$ and $yz = y$ imply $xz = x$.
(Proof: $xz = (xy)z = x(yz) = xy = x$.)

In this partial order, the product xy is the *meet* $x \wedge y$ of x and y , that is, it's the largest element z such that $z \preceq x$ and $z \preceq y$. Likewise, $x + y + xy$ is the *join* $x \vee y$ of x and y , that is, it's the smallest element z such that $x \preceq z$ and $y \preceq z$. A partial order that has meets and

joins of pairs of elements is called a *lattice*. Not all lattices have the distributive properties where meet and join distribute over each other

$$(x \vee y) \wedge z = (x \wedge z) \vee (y \wedge z) \text{ and } (x \wedge y) \vee z = (x \vee z) \wedge (y \vee z)$$

but Boolean rings do, so Boolean rings are examples of *distributive lattices*

A *minimal* element of a Boolean ring is a nonzero element such that there is no smaller nonzero element. Every element of a finite Boolean ring is a sum of the minimal elements less than or equal to it. Since there are no elements less than 0, 0 has to be treated as the empty sum.

Theorem 3.14. If R is a finite Boolean ring, then $R \cong \mathcal{P}(S)$ where

$$S = \{x \in R \mid x \text{ is minimal}\}.$$

Exercise 35. Prove the preceding theorem. Hint: see the section above 3.3.2 on factoring Boolean rings. Induction may help.

3.4 The field of rational numbers, fields of fractions

Suppose that we already have constructed the integral domain of integers \mathbf{Z} , but for some reason do not have the field of rational numbers \mathbf{Q} . Then we could construct \mathbf{Q} from \mathbf{Z} since each rational number can be named by a pair of integers. We'll do that. The steps we use only depend on \mathbf{Z} being an integral domain. That means that the construction we use can also be used to create a *field of fractions* F from any integral domain R . In the following, think of the integral domain R as being \mathbf{Z} and the field F as being \mathbf{Q} .

An equivalence relation on pairs of integers. First of all, a rational number $\frac{m}{n}$ can be named by a pair of integers (m, n) where the second integer n does not equal 0. But different pairs (m, n) and (k, l) can name the same integer $\frac{m}{n} = \frac{k}{l}$ if $ml = nk$. That suggests if we want to create rational numbers from integers, we'll need an equivalence relation on pairs of elements of the integral domain R .

We'll start with the set $R \times R_{\neq 0}$ of ordered pairs (m, n) of elements of an integral domain R with $n \neq 0$. Define a relation \equiv on this set by

$$(m, n) \equiv (k, l) \quad \text{iff} \quad ml = nk.$$

You can easily verify that this relation is an equivalence relation.

Reflexivity: $(m, n) \equiv (m, n)$. That's valid since $mn = mn$.

Symmetry: $(m, n) \equiv (k, l)$ implies $(k, l) \equiv (m, n)$. That's valid since $ml = nk$ implies $kn = lm$.

Transitivity: $(m, n) \equiv (k, l)$ and $(k, l) \equiv (s, t)$ imply $(m, n) \equiv (s, t)$. We need to show that $ml = nk$ and $kt = ls$ imply $mt = ns$. Multiply the first equation by t and the second by n . Then $mlt = nkt$ and $nkt = nls$, so $mlt = nls$. But R is an integral domain, so cancellation is valid when $l \neq 0$, so $mt = ns$.

Thus, \equiv is an equivalence relation on $R \times R_{\neq 0}$. Let F be the quotient set F_{\equiv} , and denote an element $[(m, n)]$ of F by $\frac{m}{n}$.

So far, we've got the underlying set for our proposed field F , but we don't have the operations for a field. Before we define them (and show they're well-defined), let's verify that the function $R \rightarrow R \times R_{\neq 0} \rightarrow F$ which sends an element m of R first to $(m, 1)$ then to $\frac{m}{1}$ is a one-to-one function. Suppose that $\frac{m}{1} = \frac{n}{1}$. That means $m1 = 1n$, so $m = n$. Thus we may interpret $R \rightarrow F$ as making R a subset of F by identifying m with $\frac{m}{1}$.

Addition on F . We'd like to define the sum

$$\frac{m}{n} + \frac{k}{l} \quad \text{as} \quad \frac{ml + nk}{nl},$$

but as our fractions are really equivalence classes, we need to show that's well defined. In detail, we need to show that

$$\frac{m}{n} = \frac{m'}{n'} \quad \text{and} \quad \frac{k}{l} = \frac{k'}{l'} \quad \text{imply} \quad \frac{ml + nk}{nl} = \frac{m'l' + n'k'}{n'l'}.$$

That reduces to showing that

$$mn' = nm' \quad \text{and} \quad kl' = lk' \quad \text{imply} \quad (ml + nk)n'l' = nl(m'l' + n'k').$$

But that can be shown by multiplying the first equation by ll' , the second by nn' and adding the two resulting equations. Thus, this addition on F is well-defined.

Multiplication on F . We'd like to define the product

$$\frac{m}{n} \frac{k}{l} \quad \text{as} \quad \frac{mk}{nl},$$

We need to show that's well defined. You'll find that the proof is easier than the one above for addition.

Next, we need to verify that with these definitions F satisfies the field axioms. A proof is needed for each field axiom.

Commutativity of addition. $\frac{m}{n} + \frac{k}{l} = \frac{k}{l} + \frac{m}{n}$. That's easily verified since $\frac{ml + nk}{nl} = \frac{kn + lm}{ln}$. (That depends on commutativity of addition and multiplication in R .)

Commutativity of multiplication. $\frac{m}{n} \frac{k}{l} = \frac{k}{l} \frac{m}{n}$. That's easily verified since $\frac{mk}{nl} = \frac{km}{ln}$.

Associativity of addition. You can easily show it, but it's a big mess.

Associativity of multiplication. Pretty easy.

Additive identity. $\frac{0}{1} + \frac{k}{l} = \frac{k}{l}$. Easy.

Multiplicative identity $\frac{1}{1} \frac{k}{l} = \frac{k}{l}$. Easy.

Negation. $\frac{m}{n} + \frac{-m}{n} = \frac{0}{1}$. Pretty easy.

Reciprocation. For $\frac{m}{n} \neq \frac{0}{1}$, $\frac{m}{n} \frac{n}{m} = \frac{1}{1}$. Pretty easy.

Multiplication distributes over addition. Easy but messy.

$0 \neq 1$. We need to show that $\frac{0}{1} \neq \frac{1}{1}$ in F . But that's the same as $0 \cdot 1 \neq 1 \cdot 1$ in the integral domain R , and part of the definition of integral domain requires $0 \neq 1$.

Thus, F is a field.

Exercise 36. Select four of the axioms above, and prove them. As always, your proofs should include justifications.

We'll summarize this result as a theorem.

Theorem 3.15. An integral domain R is a subring of a field F , called the *field of fractions*, where each element of F can be represented as $\frac{m}{n}$ where m and n are elements of R and $n \neq 0$.

This gives us another proof that the characteristic of an integral domain is either 0 or a prime number since it has the same characteristic of its field of fractions.

Examples 3.16. The primary example of this is the construction of \mathbf{Q} from \mathbf{Z} .

For another example, take the Gaussian integers $\mathbf{Z}[i]$ for the integral domain R . Then the field of fractions F is the field $\mathbf{Q}(i)$. The elements of $\mathbf{Q}(i)$ are of the form $x + yi$ where x and y are rational numbers.

Yet for another example, take the polynomial ring $F[x]$ with coefficients in a field F . It's an integral domain, and its field of fractions is the rational function field $F(x)$ with coefficients in F .

Stopping short of inverting all elements. Sometimes you may want to create reciprocals for some elements of an integral domain, but not for all elements. This can be done by a minor modification of the above process. Suppose, for instance, that you want to extend \mathbf{Z} to include the reciprocal of 2 but not of any other prime number. That would lead to the *domain of dyadic rationals* $\mathbf{Z}[\frac{1}{2}]$ where the denominators are powers of 2.

On the other hand, if you want to extend \mathbf{Z} to include the reciprocals of all the primes except 2, just include odd denominators. This is called localizing \mathbf{Z} at 2.

These other constructions are useful in mathematics, but we won't use them ourselves.

3.5 Categories and the category of rings

Categories are higher-order algebraic structures. We'll look at the category of rings in which the objects of the category are all the rings. The purpose of a category is to study the interrelations of its objects, and to do that the category includes morphisms between the objects. In the case of the category of rings, the morphisms are the ring homomorphisms.

We'll start with the formal definition of categories. We'll use the category of rings both to illustrate categorical concepts and to study rings. Category theory was developed by Eilenberg and Mac Lane in the 1940s.

3.5.1 The formal definition of categories

Unlike fields, rings, and groups, we won't require that categories build on sets. In a category the collection of all its objects won't be a set because the collection is larger than any set. That's not a problem since theories don't have to be built on set theory. Indeed, set theory itself is not built on set theory.

Definition 3.17. A *category* \mathcal{C} consists of

1. *objects* often denoted with uppercase letters, and
2. *morphisms* (also called *maps* or *arrows*) often denoted with lowercase letters.
3. Each morphism f has a *domain* which is an object and a *codomain* which is also an object. If the domain of f is A and the codomain is B , then we write $f : A \rightarrow B$ or $A \xrightarrow{f} B$. The collection of all morphisms from A to B is denoted $\text{Hom}(A, B)$.
4. For each object A there is a morphism $1_A : A \rightarrow A$ called the *identity morphism* on A . (When A can be determined by context, its denoted simply 1.)
5. Given two morphisms $A \xrightarrow{f} B$ and $B \xrightarrow{g} C$ where the codomain of one is the same as the domain of the other there is another morphism $A \xrightarrow{g \circ f} C$ called the *composition* of the two morphisms. This composition is illustrated by the commutative diagram

$$\begin{array}{ccc}
 A & \xrightarrow{f} & B \\
 & \searrow^{g \circ f} & \downarrow g \\
 & & C
 \end{array}$$

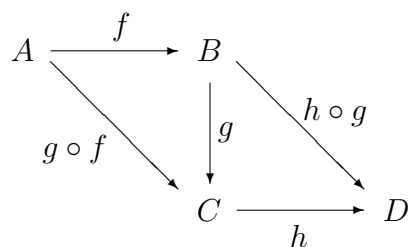
(Sometimes $g \circ f$ is denoted fg .)

A diagram of objects and morphisms in a category is said to *commute*, or be a *commutative diagram* if any two paths of morphisms (in the direction of the arrows) between any two objects yield equal compositions.

6. For all $A \xrightarrow{f} B$, $f \circ 1_A = f$ and $1_B \circ f = f$. These compositions are illustrated by the two commutative diagrams

$$\begin{array}{ccc}
 A & & A \\
 \downarrow 1_A & \searrow^{f \circ 1_A} & \downarrow f \\
 A & \xrightarrow{f} & B
 \end{array}
 \qquad
 \begin{array}{ccc}
 A & \xrightarrow{f} & B \\
 \searrow^{1_B \circ f} & & \downarrow 1_B \\
 & & B
 \end{array}$$

7. For all $A \xrightarrow{f} B$, $B \xrightarrow{g} C$, and $C \xrightarrow{h} D$, $(h \circ g) \circ f = h \circ (g \circ f)$. In the diagram below, if the two triangles in the diagram each commute, then the parallelogram commutes.



Isomorphisms in a category \mathcal{C} . Although only morphisms are defined in a category, it's easy to determine which ones are isomorphisms. A morphism $f : A \rightarrow B$ is an *isomorphism* if there exists another morphism $g : B \rightarrow A$, called its *inverse*, such that $f \circ g = 1_A$ and $g \circ f = 1_B$. Indeed, the main reason identity morphisms are included in the definition of categories is to be able to define isomorphisms.

Examples 3.18 (The categories of sets, groups, rings, and fields). Although we're more interested in the category of rings right now, the category \mathcal{S} of sets is also relevant. An object in \mathcal{S} is a set, and a morphism in \mathcal{S} is a function. The domain and codomain of a morphism are just the domain and codomain of the function, and composition is composition. Isomorphisms are bijections.

The objects of the category \mathcal{G} of groups are groups, and the morphisms of \mathcal{G} are group homomorphisms.

The objects of the category \mathcal{R} of rings are rings, and the morphisms of \mathcal{R} are ring homomorphisms.

The objects of the category of fields are fields, and its morphisms are field homomorphisms, which are just ring homomorphisms. The category of fields is a subcategory of the category of rings.

In each of these other three categories—groups, rings, fields—*isomorphisms in the category* are what we have called isomorphisms.

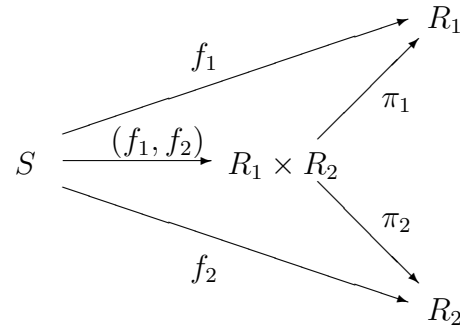
3.5.2 The category \mathcal{R} of rings

Recall that a ring homomorphism $f : A \rightarrow B$ between rings is a function that preserves addition, multiplication, and 1. The category of rings has as its objects all rings and as its morphisms all ring homomorphisms. The identity morphism 1_A on a ring is the identity homomorphism, and composition is the usual composition of homomorphisms. Thus, we have a category \mathcal{R} of rings.

If this were all there was to category theory, there wouldn't be much point to it. But by emphasizing the morphisms and deemphasizing elements in rings we can identify what's important about certain rings and certain ring constructions. We'll look at products of rings first to see what characterizes them. We'll also look at a couple of special rings, namely \mathbf{Z} and $\mathbf{Z}[x]$, for characterizing properties of them. We'll also see how to characterize monomorphisms.

The universal property of products. Recall that the product $R_1 \times R_2$ of two rings is consists of ordered pairs (x_1, x_2) with $x_1 \in R_1$ and $x_2 \in R_2$, and the ring operations for $R_1 \times R_2$ are performed coordinatewise. Furthermore, we have the projection ring homomorphisms $R_1 \times R_2 \xrightarrow{\pi_1} R_1$ and $R_1 \times R_2 \xrightarrow{\pi_2} R_2$ which pick out the two coordinates.

This product has the universal property that for each ring S and ring homomorphisms $S \xrightarrow{f_1} R_1$ and $S \xrightarrow{f_2} R_2$, there exists a unique ring homomorphism $S \rightarrow R_1 \times R_2$, which we will denote (f_1, f_2) , such that $f_1 = \pi_1 \circ (f_1, f_2)$ and $f_2 = \pi_2 \circ (f_1, f_2)$, as illustrated by the diagram below.



In fact, the product is characterized by this universal property in the sense that if another ring R has this universal property, then there is a ring isomorphism $R \rightarrow R_1 \times R_2$. In more detail, if $R \xrightarrow{p_1} R_1$ and $R \xrightarrow{p_2} R_2$ have this product property (namely, that for each ring S and ring homomorphisms $S \xrightarrow{f_1} R_1$ and $S \xrightarrow{f_2} R_2$, there exists a unique ring homomorphism $S \xrightarrow{f} R$ such that $f_1 = p_1 \circ f$ and $f_2 = p_2 \circ f$), then there exists a unique ring isomorphism $R \xrightarrow{h} R_1 \times R_2$ such that $\pi_1 \circ h = p_1$ and $\pi_2 \circ h = p_2$.

Although this characterization of products was described for the category of rings, it is the definition for the product of two objects in any category. A product $R_1 \times R_2$ is characterized by the property that a morphism to the product correspond to a pair of morphisms to the factors. The product of two sets in the category \mathcal{S} of sets has this same universal property as does the product of two groups in the category \mathcal{G} of groups. There are, however, no products in the category of fields.

\mathbf{Z} is the initial object in the category of rings. We can also use category theory to pin down what's so special about the ring \mathbf{Z} . It has the property that given any ring R , there is a unique ring homomorphism $\mathbf{Z} \xrightarrow{f} R$, and it's defined by $f(n) = n$. An object in a category with that property is called the *initial object* in the category. Any two initial objects in a category are isomorphic.

The trivial ring is the final object in the category of rings. Dual to the initial object is a final object, which in the category of rings is the *trivial* or *degenerate* ring 0 . This ring has only one element in which $0 = 1$. In fact, it's the only ring in which $0 = 1$ (since $0 = 0x = 1x = x$).

The *final object* in a category has the property that there's a unique morphism to it from each object in the category. The trivial ring has that property in the category of rings.

Exercise 37. Determine the initial object and the final object in the category \mathcal{S} of sets.

The universal property of the polynomial ring $\mathbf{Z}[x]$. Given any ring R and any element $a \in R$, there is a unique ring homomorphism $\mathbf{Z}[x] \rightarrow R$ that maps x to a . This homomorphism is just evaluation at a , and a polynomial $f(x)$ is mapped to the element $f(a)$ in R .

3.5.3 Monomorphisms and epimorphisms in a category

Although we defined a monomorphism $f : A \rightarrow B$ as a one-to-one homomorphism, we can characterize monomorphisms entirely in terms of category theory.

Definition 3.19. A morphism $f : A \rightarrow B$ is *monic*, or a *monomorphism*, when if g and h are any two morphisms from any another object C to A such that $f \circ g = f \circ h$, then $g = h$.

$$C \begin{array}{c} \xrightarrow{g} \\ \xrightarrow{h} \end{array} A \xrightarrow{f} B$$

A monomorphism in the category \mathcal{S} of sets is an injection.

This definition agrees with our previous definition for ring monomorphism in terms of elements, and one way to see the correspondence is to let C be $\mathbf{Z}[x]$. Likewise, a monomorphism in the category \mathcal{G} of groups agrees with our previous definition of group monomorphism.

Epimorphisms. The concept of epimorphism is dual to that of monomorphism. If we change the direction of all the arrows in the definition of monomorphism, we'll get the definition of epimorphism.

Definition 3.20. A morphism $f : A \rightarrow B$ is *epic*, or an *epimorphism*, when if g and h are any two morphisms from B to any another object C such that $g \circ f = h \circ f$, then $g = h$.

$$A \xrightarrow{f} B \begin{array}{c} \xrightarrow{g} \\ \xrightarrow{h} \end{array} C$$

In the category \mathcal{S} of sets, an epimorphism is a surjection. Likewise, it turns out that in the category \mathcal{G} of groups, an epimorphism is a surjection.

In the category \mathcal{R} of rings, it's easy enough to show that if f is a surjective ring homomorphism, then f is an epimorphism, but there are other epimorphisms that aren't surjections.

Example 3.21. Consider the inclusion function $\iota : \mathbf{Z} \rightarrow \mathbf{Q}$. We'll show that it's an epimorphism.

Let g and h be any two morphisms from \mathbf{Q} to any another ring C such that $g \circ \iota = h \circ \iota$. Then $g(n) = h(n)$ for any integer n . Let $\frac{m}{n}$ be a rational number with $n \neq 0$. Then $g(m) = h(m)$ and $g(n) = h(n)$. So,

$$g\left(\frac{m}{n}\right)g(n) = g\left(\frac{m}{n}n\right) = g(m) = h(m) = h\left(\frac{m}{n}n\right) = h\left(\frac{m}{n}\right)h(n) = h\left(\frac{m}{n}\right)g(n).$$

Since $n \neq 0$, therefore $g(n) \neq 0$ as well. Cancel the $g(n)$ at the ends of the continued equation to conclude $g\left(\frac{m}{n}\right) = h\left(\frac{m}{n}\right)$. Thus, $g = h$.

Therefore, the ring homomorphism $\iota : \mathbf{Z} \rightarrow \mathbf{Q}$ is an epimorphism in \mathcal{R} , the category of rings. It is also a monomorphism. But it is not an isomorphism.

In many categories, if a morphism is both monic and epic, then it's also an isomorphism. That's true in the category \mathcal{S} of sets and in the category \mathcal{G} of groups, but not in the category \mathcal{R} of rings. This example shows that \mathcal{R} is a somewhat unusual category.

3.6 Kernels, ideals, and quotient rings

These three concepts are closely related. For a ring homomorphism $f : R \rightarrow S$, the inverse image of 0 is a subset of R called the kernel of f and denoted $\text{Ker } f$. It can't be just any subset, as we'll see, since it's closed under addition and multiplication by elements of R . A subset with those properties we'll call an ideal of R . Every ideal I of R is the kernel of some ring homomorphism $f : R \rightarrow S$. We'll use an ideal I of a ring R to define a quotient ring R/I and a projection $\gamma : R \rightarrow R/I$. These projections will be generalizations of the projections $\mathbf{Z} \rightarrow \mathbf{Z}_n$ that we studied earlier.

3.6.1 Kernels of ring homomorphisms

Definition 3.22. Let $f : R \rightarrow S$ be a ring homomorphism. Those elements of R that are sent to 0 in S form the *kernel* of f .

$$\text{Ker } f = f^{-1}(0) = \{x \in R \mid f(x) = 0\}.$$

We'll look at properties of this kernel and see what it tells us about the function f .

Example 3.23. It's a good idea to have in mind an example or two whenever a new concept is defined. The definition of the kernel of a ring homomorphism is given above, and a good example for it is the ring homomorphism $f : \mathbf{Z} \rightarrow \mathbf{Z}_n$ where n is a fixed integer. That's an especially good example we can use it throughout this discussion of rings, ideals, and quotient rings.

For that $f : \mathbf{Z} \rightarrow \mathbf{Z}_n$, an element $x \in \mathbf{Z}$ is in $\text{Ker } f$ if it is sent to $[0]_n$, the 0 element in the ring \mathbf{Z}_n , that is, if $[x]_n = [0]_n$, or, more simply, if $n \mid x$. Therefore, the kernel of f consists of the multiples of n . A standard notation for the multiples of an integer n is $n\mathbf{Z}$. Thus, for this function f , $\text{Ker } f = n\mathbf{Z}$.

Kernels aren't just any subsets of R ; they have some special properties. We have, of course, $0 \in \text{Ker } f$, since $f(0) = 0$. Also, if x and y are both in $\text{Ker } f$, then $f(x + y) = f(x) + f(y) = 0 + 0 = 0$, so their sum $x + y$ is also in $\text{Ker } f$. Furthermore, if $x \in \text{Ker } f$ and y is any element of R , then $f(xy) = f(x)f(y) = 0f(y) = 0$, so $xy \in \text{Ker } f$, and likewise $yx \in \text{Ker } f$.

Besides telling us what elements are sent to 0 by f , the kernel of f also tells us when two elements are sent to the same element. Since $f(x) = f(y)$ if and only if $f(x - y) = 0$, therefore, f will send x and y to the same element of S if and only if $x - y \in \text{Ker } f$.

3.6.2 Ideals of a ring

The properties of kernels of homomorphisms that we just found we'll use to define ideals of rings. Historically, ideals had a different purpose, but we'll get to that purpose later. The word "ideal" is short for ideal number or ideal element.

Definition 3.24. An *ideal* I of a ring R is a subset that (1) includes 0, (2) is closed under addition, and (3) is closed under multiplication by elements of R . We can summarize these requirements symbolically by $0 \in I$, $I + I \subseteq I$, $RI \subseteq I$, and $IR \subseteq I$.

Both of the last two requirements, $RI \subseteq I$ and $IR \subseteq I$, are needed when R is a non-commutative ring. Most of the time we'll be dealing with commutative rings so one will do.

Note that $\{0\}$ is always an ideal in a ring R . It's called the *trivial* ideal. We'll usually just denote it 0. Also, the entire ring R is an ideal, but not a proper ideal. A *proper ideal* is any ideal $I \neq R$.

Theorem 3.25. The intersection of ideals is an ideal.

Proof. Here's the proof for two ideals I_1 and I_2 of a ring R . This proof can be generalized to any number, including an infinite number, of ideals.

We need to show that $I_1 \cap I_2$ (1) includes 0, (2) is closed under addition, and (3) is closed under multiplication by elements of R .

First, since $0 \in I_1$ and $0 \in I_2$, therefore $0 \in I_1 \cap I_2$.

Second, given two elements $x, y \in I_1 \cap I_2$, to show $x + y \in I_1 \cap I_2$. Since $x, y \in I_1 \cap I_2$, therefore $x, y \in I_1$ and $x, y \in I_2$. Therefore $x + y \in I_1$ and $x + y \in I_2$, and so $x + y \in I_1 \cap I_2$.

Third, given $x \in I_1 \cap I_2$ and $y \in R$, to show $xy \in I_1 \cap I_2$. Since $x \in I_1 \cap I_2$, therefore $x \in I_1$ and $x \in I_2$. Therefore, $xy \in I_1$ and $xy \in I_2$, and so $xy \in I_1 \cap I_2$. Q.E.D.

Principal ideals and ideals generated by a set. The simplest examples of ideals are what are called principal ideals. Let a be an element of a commutative ring R . The set of all multiples of a ,

$$(a) = \{xa \mid x \in R\},$$

is an ideal of R , as you can easily check. These ideals are called *principal ideals* because they are generated by one element. An alternate notation for the principal ideal generated by the element a is Ra or aR .

Note that (0), the ideal generated by 0, is just the 0 ideal, while (1), the ideal generated by 1, is all of R .

Sometimes it takes more than one element to generate an ideal. Let A be a subset of a commutative ring R . The smallest ideal that contains A is called the *ideal generated by A*. It must contain all linear combinations of elements of A since an ideal is closed under addition and closed under multiplication by elements of R , but that's enough. Usually, we're only interested in generating an ideal from a finite number of elements $A = \{a_1, a_2, \dots, a_k\}$. Then the ideal generated by A is

$$(a_1, a_2, \dots, a_k) = \{x_1a_1 + \dots + x_ka_k \mid \text{each } x_i \in R\}.$$

An example of an ideal generated by two elements but not principal (not by one element) is $(5, x^2)$ in $\mathbf{Z}[k]$, the polynomial ring with integral coefficients.

Exercise 38. As you know, if $n \in \mathbf{Z}$, then $n\mathbf{Z}$, also written (n) , is an ideal of the ring \mathbf{Z} . Consider the two ideals $I = 6\mathbf{Z}$ and $J = 10\mathbf{Z}$ of the \mathbf{Z} .

- (a). Determine their intersection $I \cap J$ as a principal ideal of \mathbf{Z} .
- (b). Prove that the union $I \cup J$ is not an ideal of \mathbf{Z} .

3.6.3 Quotient rings, R/I

As mentioned above the kernel of a ring homomorphism f tells us when two elements are sent to the same element: $f(x) = f(y)$ if and only if $x - y \in \text{Ker } f$. We can use $\text{Ker } f$ to construct a “quotient ring” $R/\text{Ker } f$ by identifying two elements x and y in R if their difference lies in $\text{Ker } f$. In fact, we can do this not just for kernels of homomorphisms, but for any ideal I . That is, we can use an ideal I of R to determine when two elements x and y are to be identified, $x \equiv y$, and we’ll end up with a ring R/I . The identification is called a congruence. This concept of congruence generalizes congruence modulo n on \mathbf{Z} .

Definition 3.26. A *congruence* \equiv on a ring R is an equivalence relation such that for all $x, x', y, y' \in R$,

$$x \equiv x' \text{ and } y \equiv y' \text{ imply } x + y \equiv x' + y' \text{ and } xy \equiv x'y'.$$

Since we’re dealing with rings with 1, we’ll usually insist that $0 \neq 1$. The equivalence classes for a congruence are called *congruence classes*.

Theorem 3.27. If \equiv is a congruence on a ring R , then the quotient set R/\equiv , that is, the set of congruence classes, is a ring where addition is defined by $[x] + [y] = [x + y]$ and multiplication by $[x][y] = [xy]$.

Proof. First we need to show that the proposed definitions are actually well defined. That is, if a different representative x' is chosen from the congruence class $[x]$ and y' from $[y]$, then the same classes $[x' + y']$ and $[x'y']$ result. That is

$$[x] = [x'] \text{ and } [y] = [y'] \text{ imply } [x + y] = [x' + y'] \text{ and } [xy] = [x'y'].$$

That’s the same as the requirements met in the definition of congruence (which explains why they are in the definition).

Also, each of the axioms for a ring need to be verified, but they’re all automatic. Here’s commutativity of addition, for example.

$$[x] + [y] = [x + y] = [y + x] = [y] + [x].$$

We could say that the quotient ring inherits the properties from the ring.

Q.E.D.

In the next theorem we’ll see that an ideal I determines a congruence. We’ll write the congruence $x \equiv y \pmod{I}$ rather than just $x \equiv y$ when we want to emphasize the role of I . The congruence classes may be written $[x]$ or $[x]_I$, or $x + I$. The last notation is a good one since $[x] = \{x + y \mid y \in I\}$.

Theorem 3.28 (Congruence modulo an ideal). Let I be an ideal of a ring R . A congruence, called *congruence modulo I* , is defined by

$$x \equiv y \pmod{I} \text{ if and only if } x - y \in I.$$

The quotient ring, R/\equiv , is denoted R/I .

Proof. First, we need to show that it's an equivalence relation.

Reflexivity. $x \equiv x \pmod{I}$. That's okay since $x - x = 0 \in I$.

Symmetry. $x \equiv y \pmod{I}$ implies $y \equiv x \pmod{I}$. That's okay because if $x - y \in I$, then $y - x = -(x - y) \in I$.

Transitivity. $x \equiv y \pmod{I}$ and $y \equiv z \pmod{I}$ imply $x \equiv z \pmod{I}$. That's okay, too. If $x - y \in I$ and $y - z \in I$, then so is their sum $x - z \in I$.

Thus, it's an equivalence relation. Next to show that

$$x \equiv x' \pmod{I} \text{ and } y \equiv y' \pmod{I} \text{ imply } x + y \equiv x' + y' \pmod{I} \text{ and } xy \equiv x'y' \pmod{I}.$$

That requirement reduces to the statement

$$x - x' \in I \text{ and } y - y' \in I \text{ imply } (x + y) - (x' + y') \in I \text{ and } (xy - x'y') \in I,$$

which, you can check, follow from the definition of ideal.

Q.E.D.

Exercise 39. Prove the last statement above: if $x - x' \in I$ and $y - y' \in I$, then $(x + y) - (x' + y') \in I$ and $(xy - x'y') \in I$.

Example 3.29 (Cyclic rings). As we saw above, $I = n\mathbf{Z}$ is an ideal of \mathbf{Z} . The congruence defined here is the same one we had before. Thus, $x \equiv y \pmod{I}$ means $x \equiv y \pmod{n}$. The quotient ring is $\mathbf{Z}/n\mathbf{Z}$, which we have studied before and denoted \mathbf{Z}_n for short.

Comment 3.30. The ring structure on the quotient R/I was defined from the ring structure on R , so the projection $\gamma : R \rightarrow R/I$ is a ring homomorphism. This ring R/I is called a *quotient ring* of R . (It is also sometimes called a factor ring, but that term should be restricted to the case when R factors as a product of rings, one of which is R/I . An example of that is the Chinese remainder theorem.)

Examples 3.31 (Quadratic field extensions.). We've looked at $\mathbf{Q}(\sqrt{2})$, $\mathbf{C} = \mathbf{R}(i)$, and other quadratic field extensions. We can interpret them as quotient rings.

Let's take $\mathbf{Q}(\sqrt{2})$ first. Consider the ring $R = \mathbf{Q}[x]$ of polynomials with rational coefficients. An ideal in R is the principal ideal $I = (x^2 - 2)$ generated by the polynomial $x^2 - 2$. In the quotient ring $R/I = \mathbf{Q}[x]/(x^2 - 2)$, we have $x^2 - 2 \equiv 0 \pmod{I}$, that is, $x^2 \equiv 2 \pmod{I}$, so in R/I , we find that 2 does have a square root, namely x . Since in R/I every polynomial $a_n x^n + \cdots + a_1 x + a_0$ is congruent to a polynomial of degree 1 (because $x^2 \equiv 2$), but no two linear polynomials are congruent mod I (because $a_1 x + a_0 \equiv b_1 x + b_0 \pmod{I}$ implies $(a_1 - b_1)x + (a_0 - b_0) \in I$ so $a_1 = b_1$ and $a_0 = b_0$), therefore every element in R/I is uniquely represented as a linear polynomial $a_1 x + a_0$. If we denote x by the symbol $\sqrt{2}$, then we find $\mathbf{Q}[x]/(x^2 - 2)$ is the same field as $\mathbf{Q}(\sqrt{2})$ that we described before.

Likewise, $\mathbf{R}[x]/(x^2 + 1)$ is \mathbf{C} .

We'll find this construction of new rings as quotient rings is very useful, especially when we take quotients rings of polynomial rings like we did here.

The image of a ring homomorphism is isomorphic to the ring modulo its kernel.

Let $f : R \rightarrow S$ be a ring homomorphism. The image of f , denoted $f(R)$, is the set

$$f(R) = \{f(x) \in S \mid x \in R\}.$$

It is a subring of S , as you can easily verify. You can also show the following isomorphism theorem, called the first isomorphism for rings.

Theorem 3.32. If $f : R \rightarrow S$ is a ring homomorphism then the quotient ring $R/\text{Ker } f$ is isomorphic to the image ring $f(R)$, the isomorphism being given by

$$\begin{aligned} R/\text{Ker } f &\rightarrow f(R) \\ x + \text{Ker } f &\mapsto f(x) \end{aligned}$$

Exercise 40. Prove the preceding theorem.

(a). First show that the assignment $x + \text{Ker } f$ to $f(x)$ is well defined. That means that if $x + \text{Ker } f = x' + \text{Ker } f$, then $f(x) = f(x')$. Call that function $\phi(x)$.

(b). Show that assignment is a ring homomorphism. Show (1) $\phi(1) = 1$, (2) $\phi(x + y) = \phi(x) + \phi(y)$, and (3) $\phi(xy) = \phi(x)\phi(y)$.

This gives us two ways to look at the image, either as a quotient ring of the domain R or as a subring of the codomain S .

Furthermore, we can now treat a ring homomorphism $f : R \rightarrow S$ as a composition of three ring homomorphisms.

$$R \rightarrow R/\text{Ker } f \cong f(R) \rightarrow S$$

The first is the projection from R onto its quotient ring $R/\text{Ker } f$, the second is the isomorphism $R/\text{Ker } f \cong f(R)$, and the third is the inclusion of the image $f(R)$ as a subring of S .

3.6.4 Prime and maximal ideals

Sometimes it occurs that R/I is not just a ring, but either an integral domain or even a field. Those results occur when the ideal I is a prime ideal or a maximal ideal, respectively, as we'll define now.

Definition 3.33. An ideal I in a commutative ring R is said to be a *prime ideal* if R/I is an integral domain. Equivalently, I is a prime ideal if (1) $I \neq R$, and (2) $\forall x, y \in R$, if $xy \in I$, then either $x \in I$ or $y \in I$. An ideal I is said to be *maximal* if it's a proper ideal, but it is not contained in any larger proper ideal.

Exercise 41. Prove that R/I is an integral domain if and only if R/I satisfies both conditions (1) $I \neq R$, and (2) $\forall x, y \in R$, if $xy \in I$, then either $x \in I$ or $y \in I$.

Example 3.34. The ideals of \mathbf{Z} that are prime are those of the form $p\mathbf{Z}$ where p is a prime number, and the 0 ideal. In fact, $p\mathbf{Z}$ are maximal ideals, but 0 is not maximal.

In a field F there is only one proper ideal, namely 0.

In an integral domain, the 0 ideal is a prime ideal, and conversely, if 0 is an ideal in a commutative ring, then the ring is an integral domain.

Theorem 3.35. Every maximal ideal is prime.

Proof. Let I be a maximal ideal of a commutative ring R , and let $xy \in I$. Suppose $x \notin I$. Then $xR + I = \{xu + v \mid u \in R, v \in I\}$ is an ideal containing I . Since I is an maximal ideal, therefore $xR + I$ is not a proper ideal but all of R . Therefore $1 = xu + v$ for some $u \in R, v \in I$. Hence $y = yxu + yv \in Iu + I = I$. Thus, I satisfies the conditions to be a prime ideal. Q.E.D.

We won't show it right now, but we'll prove later Krull's theorem which says that every ideal is contained in a maximal ideal. We'll need to discuss the axiom of choice and Zorn's lemma before we can prove it.

Theorem 3.36 (Maximal ideal theorem). Let I be an ideal of a commutative ring R . Then I is a maximal ideal if and only if R/I is a field.

Proof. We'll use the notation $[x]$ for $x + I$ to stress that we're thinking of it as an element of R/I .

Suppose that I is a maximal ideal, and let $[x]$ be any nonzero element of R/I , that is $x \notin I$. As in the last proof, $xR + I = R$. Therefore $1 = xu + v$ for some $u \in R, v \in I$. Then, in R/I we have $[1] = [x][u] + [v] = [x][u] + [0] = [x][u]$. Therefore $[x]$ has a reciprocal, and R/I is a field.

Now suppose that R/I is a field. Let $x \notin I$. We'll show that $xR + I = R$ which will show that I is a maximal ideal. In R/I , $[x] \neq [0]$, so $[x]$ has an inverse $[y]$, $[x][y] = [1]$, so $1 - xy \in I$, so $1 \in xR + I$, hence $R = xR + I$. Q.E.D.

3.7 Krull's theorem

We'd like to prove Krull's theorem that every ideal in a commutative ring is contained in a maximal ideal, but in order to do that in general we'll need something called Zorn's lemma. It's a statement that's logically equivalent to the better known axiom of choice.

See section A.4 in the appendix for a review of the axiom of choice and Zorn's lemma.

Theorem 3.37 (Krull). Let I be a proper ideal of a commutative ring R . Then there is a maximal ideal J such that $I \subseteq J$.

Proof. Consider the collection \mathcal{M} of proper ideals of R that contain I . Note that \mathcal{M} is nonempty since $I \in \mathcal{M}$.

We'll show that every chain \mathcal{C} in \mathcal{M} has an upper bound in \mathcal{M} . Let $B = \bigcup_{A \in \mathcal{C}} A$. Certainly B is an upper bound for \mathcal{C} since B is just the union of elements of \mathcal{C} .

We still have to show B is an ideal, which requires $RB \subseteq B$ and $B + B \subseteq B$. For the first, $RB = R \left(\bigcup_{A \in \mathcal{C}} A \right) = \bigcup_{A \in \mathcal{C}} RA = \bigcup_{A \in \mathcal{C}} A = B$. Now let $x, y \in B$. Then $x \in A_1$ for some $A_1 \in \mathcal{C}$ and $y \in A_2$ for some $A_2 \in \mathcal{C}$. But \mathcal{C} is a chain, so either $A_1 \subseteq A_2$ or $A_2 \subseteq A_1$. In the first case, $x, y \in A_2$, so $x + y \in A_2 \subseteq B$, and in the second $x, y \in A_1$, so $x + y \in A_1 \subseteq B$. Thus, $B + B \subseteq B$.

Now we can apply Zorn's lemma. It implies \mathcal{M} has a maximal element J . Clearly, $I \subseteq J$, and J is a proper ideal of R , but there are no larger proper ideals of R that contain J , so J is a maximal ideal. Q.E.D.

Note how we have not actually found J . There may be many different maximal ideals that contain I , and one was selected by a choice function, but we don't even know what the choice function is so we can't even determine J in principle.

It's actually the case that Krull's theorem is logically equivalent to the Axiom of Choice. That is, if Krull's theorem is taken as an axiom, then the Axiom of Choice can be proved from it.

There are many other applications of Zorn's lemma. For instance, you can prove that every vector space has a basis, even when the vector space is infinite dimensional.

3.8 Unique factorization domains, principal ideal domains, and Euclidean domains

Not every integral domain is as nice as the ring of integers. The ring of integers has three nice properties. One is unique factorization—every integer is uniquely a product of prime numbers. A second is that every ideal is a principal ideal. A third is that there is a division algorithm that is the basis of the Euclidean algorithm.

There aren't many rings that have all these properties, and some rings have none of them. We'll investigate these properties and their interrelations.

We'll use these three properties to define three special kinds of integral domains: unique factorization domains (UFDs), principal ideal domains (PIDs), and Euclidean domains (EDs). When we do we'll find every Euclidean domain is a principal ideal domain, every principal ideal domain is a unique factorization domain, every unique factorization domain is an integral domain; and every integral domain is a ring.

$$\text{EDs} \subset \text{PIDs} \subset \text{UFDs} \subset \text{Integral domains} \subset \text{Commutative rings}$$

3.8.1 Divisibility in an integral domain

We'll borrow the concepts of divisibility and greatest common divisor from \mathbf{Z} and apply them to integral domains. We'll separate the concept of prime number in \mathbf{Z} into two concepts since in some of the integral domains we'll look at they're actually different.

Definition 3.38. The following definitions apply to elements of an integral domain.

- Let a and b be nonzero elements. We'll say a divides b , written $a|b$, if there exists c such that $ac = b$.
- We'll say that d is a *greatest common divisor* of a and b , if d divides both a and b , and whenever another element e divides both a and b , then e divides d .
- An element x that is not zero and not a unit is *irreducible* if whenever $x = yz$, either y or z is a unit, otherwise it is *reducible*.

- An element x that is not zero and not a unit is *prime* if whenever $x|yz$, then $x|y$ or $x|z$.

Note that we won't use the notation $d = \text{GCD}(a, b)$ when d is a greatest common divisor since there will be other greatest common divisors, that is, the greatest common divisor is only unique up to a unit. Later, when we look at principal ideal domains, we can use the notation $(c) = (a, b)$ for greatest common divisors which says the principal ideal (c) is the same as the ideal generated by a and b .

Exercise 42. Several properties of divisibility follow directly from the definition just like they do with the integral domain \mathbf{Z} . Prove the following properties from the above definitions.

- 1 divides every element.
- Each element divides itself.
- If $a|b$ then $a|bc$.
- Divisibility is transitive.
- If one element divides two other elements, then it divides both their sum and difference.
- Cancellation: When $c \neq 0$, $a|b$ if and only if $ac|bc$.

Theorem 3.39. If an element in an integral domain is prime, then it is irreducible.

Proof. Let x be prime. Suppose that $x = yz$. Then $x|yz$, so either $x|y$ or $x|z$. In the first case, $xw = y$ for some w . Therefore $xwz = yz = x$. Cancel the x to conclude $wz = 1$. Then z is a unit. Likewise, in the second case y is a unit. Therefore x is irreducible. Q.E.D.

The converse of this theorem does not hold. That is, there are integral domains where not all irreducible elements are prime. We'll see that in this next example. But then a little later, we'll see that in principal ideal domains (about to be defined), irreducible elements are prime.

Example 3.40 (a nonUFD). We'll find a number of other UFDs, but, it's important to know that not every integral domain has unique factorization. Consider the integral domain $R = \mathbf{Z}[\sqrt{10}]$. An element of it is of the form $x + y\sqrt{10}$ where x and y are integers. In this integral domain 9 can be factored in two ways.

$$9 = 3^2 = (\sqrt{10} + 1)(\sqrt{10} - 1),$$

but 3, $\sqrt{10} + 1$, and $\sqrt{10} - 1$ are all irreducible. This integral domain, and many others, are not UFDs. Although the three elements 3, $\sqrt{10} + 1$, and $\sqrt{10} - 1$ are irreducible, none divides any other, so none of them is prime, as you can see by the equation involving 9, above.

3.8.2 Unique factorization domains

Unique factorization is a property that we might expect, but it turns out it doesn't hold in every integral domain. Given any element x in a ring D , we expect that we can factor it into 'atoms,' things that can't be cut further, and that there's only one way to do that. Of course, with our experience with the integers, we know that there's a bit of difficulty in stating the uniqueness part of the claim. For one thing, the order of the factors is variable, and, for another, there are units, like 1 and -1 that can be inserted to change the formal listing of the factors. Still, these are small things that we can deal with.

Definition 3.41. An integral domain is a *unique factorization domain* (UFD) if every element in it is a product of irreducible elements and it is a product of irreducible elements in only one way apart from the order of the product and factors of units.

The ring \mathbf{Z} of integers is, of course, a unique factorization domain. An integer, such as 6 can be written in more than one way as a product of irreducible elements (primes, in the case of integers) $6 = 2 \cdot 3 = (-3) \cdot (-2)$, but the only difference is the order of the primes and the insertions of units in the factorization.

Recall that an ideal I in a commutative ring R is a prime ideal when R/I is an integral domain. Equivalently, I is a prime ideal if (1) $I \neq R$, and (2) for all $x, y \in R$, if $xy \in I$, then either $x \in I$ or $y \in I$.

Theorem 3.42. An nonzero element x in an integral domain D is prime if and only if the principal ideal (x) is a prime ideal.

Exercise 43. Prove the preceding theorem. Note that there are two things to prove in an if-and-only-if statement.

3.8.3 Principal ideal domains

A second nice property that the ring of integers has is that every ideal in \mathbf{Z} is generated by a single element. If I is an ideal in \mathbf{Z} , then the GCD of all its nonzero elements is an element of I and all other elements are multiples of this GCD. This will be our definition of a principal ideal domain (PID), and we'll show that every PID is a UFD. There are UFDs that aren't PIDs, for instance, $\mathbf{Z}[x]$, the ring of polynomials with integer coefficients is one; one nonprincipal ideal is generated by 2 and x .

Definition 3.43. An integral domain is a *principal ideal domain* (PID) if every ideal in the domain is principal, that is, generated by one element.

Besides \mathbf{Z} , other prominent PIDs are $F[x]$ where F is a field. We'll prove this in section 3.8.4 on Euclidean domains which are special kinds of PIDs.

We'll show in a couple of steps that every PID is a UFD. The first one makes a connection between greatest common divisors and ideals.

Theorem 3.44. Let D be a principal ideal domain with nonzero elements a and b . The ideal (a, b) is principal, so it is equal to (c) for some element c . Then c is a greatest common divisor of a and b .

Proof. Since $a \in (c)$, therefore $c|a$. Likewise, $c|b$. We also know that $c \in (a, b)$, so $c = xa + yb$ for some elements x and y .

To show that c is a greatest common divisor, suppose d is some other common divisor of a and b . Then $a = ud$ and $b = vd$ for some elements u and v . Now,

$$c = xa + yb = xud + yvd = (xu + yv)d.$$

Therefore, $d|c$. Thus c is a greatest common divisor of a and b .

Q.E.D.

Theorem 3.45. In a principal ideal domain, irreducible elements are prime.

Proof. Suppose that p is irreducible and $p|ab$. We'll show either $p|a$ or $p|b$. We'll do that by showing that if p doesn't divide a , then it does divide b .

Suppose p does not divide a . Then the ideal (p, a) is (1) since p is irreducible. Since $1 \in (p, a)$, $1 = xp + ya$ for some elements x and y . Therefore, $b = bxp + aby$. Since $p|ab$, therefore $p|b xp + aby$, so $p|b$.

Thus, the irreducible element p is also prime. Q.E.D.

Next, we'll use the following lemma to show that elements have factorizations in PIDs. We'll still have to show they're unique after that. The condition in the lemma is called the *ascending chain condition* (ACC) on ideals, and rings that satisfy it are called *Noetherian rings* in honor of Noether who studied such rings.

Lemma 3.46. In a principal ideal domain, there are no infinitely ascending chains of ideals. That is,

$$(a_1) \subsetneq (a_2) \subsetneq (a_3) \subsetneq \cdots$$

does not exist.

Proof. Suppose there were such an infinitely ascending chain of ideals. Then the union $I = \bigcup_{i=1}^{\infty} (a_i)$ is an ideal, as you can easily check. It must be principal, so $I = (a)$ for some element a . But a is in the union, so it's in one of the ideals (a_i) . Then

$$(a) \subseteq (a_i) \subsetneq (a_{i+1}) \subseteq (a),$$

a contradiction. Q.E.D.

There are rings, in fact UFDs, that are not Noetherian. An example is a polynomial ring with infinitely many variables such as $\mathbf{Q}[x_1, x_2, x_3, \dots]$. An infinitely ascending chain of ideals in that ring is $(x_1) \subsetneq (x_1, x_2) \subsetneq (x_1, x_2, x_3) \subsetneq \cdots$.

Theorem 3.47. In a principal ideal domain, every element that is not zero and not a unit has a factorization into irreducible elements.

Proof. Suppose that a nonzero element a_1 has no factorization into irreducible elements. We'll derive a contradiction, but we'll need an element with no factorization with an extra property. We'll get that element, denoted a_n below, as follows.

Starting with the ideal (a_1) , form any ascending chain of ideals generated by other elements with no factorizations, and extend the chain as far as possible. By the lemma, it stops somewhere, say at (a_n) .

$$(a_1) \subsetneq (a_2) \subsetneq \cdots \subsetneq (a_n).$$

We now have an element a_n which has no factorization into irreducible elements with an extra property, namely, any ideal strictly containing (a_n) is generated by an element that does have such a factorization. Now, a_n is not irreducible itself, for that would be a factorization, so $a_n = bc$ where neither b nor c is a unit. Since $b|a_n$, therefore $(a_n) \subseteq (b)$. But $(a_n) \neq (b)$, for otherwise $b = a_n d$ for some d , and then $a_n d c = bc = a_n$, so $dc = 1$ making c a unit, which it is not.

So $(a_n) \subsetneq (b)$ and likewise $(a_n) \subsetneq (c)$, therefore both b and c have factorizations, and the product of those factorizations gives a factorization for a_n , a contradiction. Q.E.D.

Theorem 3.48. Every principal ideal domain is a unique factorization domain.

Proof. The last theorem gave the existence of at least one factorization for an element a . We still have to show that there's at most one factorization.

Suppose that a has two factorizations as products of irreducible elements.

$$a = p_1 \cdots p_n = q_1 \cdots q_m$$

Since the irreducible element p_1 is prime (in a PID), p_1 divides one of the q_i 's, which we can renumber as q_1 . Then $p_1 = u_1 q_1$ where u_1 is a unit. Substitute $u_1 q_1$ for p_1 , and cancel q_1 to get the equation

$$u_1 p_2 \cdots p_n = q_2 \cdots q_m.$$

That completes the inductive step of mathematical induction on n . The base case, when $n = 1$, is left to the reader. Q.E.D.

3.8.4 Euclidean domains

The third nice property that \mathbf{Z} has is that there is a division algorithm that is the basis of the Euclidean algorithm.

Some example Euclidean domains besides \mathbf{Z} that we'll discuss in this section include the Gaussian integers $\mathbf{Z}[i]$, the Eisenstein integers $\mathbf{Z}[\omega]$ where ω is a primitive cube root of 1, and polynomial rings $F[x]$ over a field F .

For the integers, the division algorithm starts with an integer a (the dividend) and a nonzero integer b (the divisor) and delivers q (the quotient) and r (the remainder) such that

$$a = qb + r \quad \text{and} \quad 0 \leq r < b.$$

This property allowed us to construct the Euclidean algorithm for finding GCDs as well as the extended Euclidean algorithm to show that the greatest common divisor of two numbers is a linear combination of them.

There are a few other integral domains that have the same kind of division algorithm where the remainder is somehow "smaller" than the divisor, but the concept of smaller and how to find q and r differs from domain to domain.

Definition 3.49. A *Euclidean valuation* on an integral domain D is a function $v : D - 0 \rightarrow \mathbf{Z}_{\geq 0}$ that satisfies the conditions

1. for nonzero elements a and b , $v(a) \leq v(ab)$, and
2. for each element a (the dividend) and nonzero element b (the divisor), there are elements q (the quotient) and r (the remainder) such that

$$a = qb + r \quad \text{where either } r = 0 \text{ or } v(r) < v(b).$$

An integral domain that admits a Euclidean valuation is called *Euclidean domain*.

Of course, \mathbf{Z} is a Euclidean domain with the valuation being the absolute value $v(a) = |a|$.

Another class of Euclidean domains are the rings of polynomials (in one variable) with coefficients in a given field. The following theorem is essentially just long division for polynomials. We'll make it simple by making the divisor $g(x)$ a monic polynomial, that is, a polynomial whose leading coefficient is 1.

It directly follows from the division algorithm for polynomials over a field, theorem 1.55, that a field's polynomial ring is a Euclidean domain.

Corollary 3.50. The polynomial ring $F[x]$ with coefficients in a field F is a Euclidean domain where the valuation v assigns to a polynomial $f(x)$ the degree of f .

Soon we'll study polynomial rings in more detail.

There are other Euclidean domains including the Gaussian integers and the Eisenstein integers.

The Gaussian integers $\mathbf{Z}[i]$ is a Euclidean domain. The ring of Gaussian integers is $\mathbf{Z}[i] = \{a_1 + a_2i \mid a_1, a_2 \in \mathbf{Z}\}$. Its valuation function, also called the norm, is $v(a_1 + a_2i) = a_1^2 + a_2^2$, the square of the distance to the origin. In order to divide one Gaussian integer $b_1 + b_2i$ into another $a_1 + a_2i$ to get a quotient $q_1 + q_2i$ and remainder $r_1 + r_2i$, you can perform the complex division $\frac{a_1 + a_2i}{b_1 + b_2i}$ to get an exact quotient, and choose $q_1 + q_2i$ to be the closest Gaussian integer to that exact quotient. The remainder is then determined.

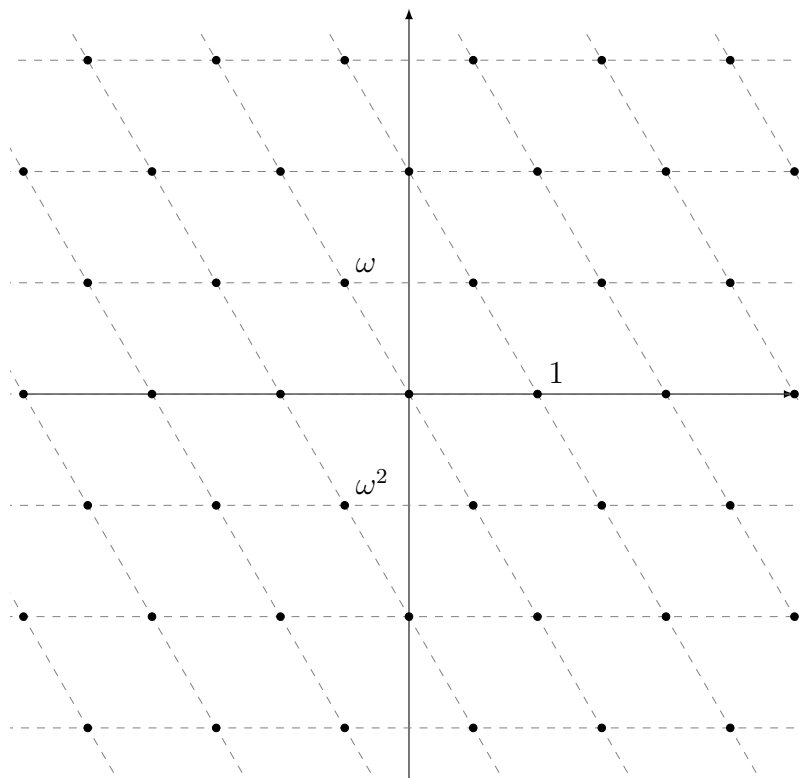


Figure 3.3: Lattice of Eisenstein integers

Eisenstein integers Whereas a basis for the Gaussian integers consists of 1 and i , a basis for the Eisenstein integers consists of 1 and ω where $\omega = \frac{1}{2}(-1 + i\sqrt{3})$ is a primitive cube root of unity. A primitive cube root of unity satisfies the equation $\frac{x^3 - 1}{x - 1} = 0$ which simplifies to $x^2 + x + 1 = 0$. The lattice of Eisenstein integers is a triangular lattice since $1 + \omega + \omega^2 = 0$. The lattice is shown in figure 3.3. The dotted lines show coordinates relative to the basis consisting of 1 and ω .

There are six units in the Eisenstein integers. They are the six sixth roots of unity: 1 itself, a primitive sixth root $\omega - 1$, a primitive cube root ω , the primitive square root -1 , and a primitive sixth root $-\omega$. They are equally spaced at 60° around the unit circle.

Like the Gaussian integers, the Eisenstein integers also are a Euclidean domain. The valuation is $v(a + b\omega) = a^2 - ab + b^2$.

The Euclidean algorithm in Euclidean domains. First, we'll show that Euclidean domains are principal ideal domains, and since PIDs are also UFDs, therefore Euclidean domains are also unique factorization domains. Then we'll look at an example of the Euclidean algorithm in a Euclidean domain other than \mathbf{Z} .

Theorem 3.51. A Euclidean domain is a principal ideal domain.

Proof. Let I be an ideal in a Euclidean domain D with valuation v . We'll show I is a principal ideal. If I is the zero ideal (0) , then it's principal of course.

Assume now that I has a nonzero element, and let $S = \{v(x) \mid 0 \neq x \in I\}$. This is a nonempty subset of the nonnegative integers, so it has a least element, and let that be $v(a)$. Thus, a is a nonzero element of I , so $(a) \subseteq I$. Let x be any other nonzero element in I . Then $v(a) \leq v(x)$. Furthermore, there are elements q and r in D such that $x = aq + r$ and either $r = 0$ or $v(r) < v(a)$. But $r = x - aq \in I$, so if $r \neq 0$, then $v(r) > v(a)$ contradicts $v(a) \leq v(r)$. Therefore, $r = 0$, and hence $x = aq$, so $a \mid x$. Therefore, $I = (a)$. Thus, D is a PID. Q.E.D.

The Euclidean algorithm works in any Euclidean domain the same way it does for integers. It will compute the greatest common divisor (up to a unit), and the extended Euclidean algorithm will construct the greatest common divisor as a linear combination of the original two elements.

Example 3.52. Let's take an example from the polynomial ring $\mathbf{Q}[x]$. Let's find the greatest common divisor of $f_1(x) = x^4 + 2x^3 - x - 2$ and $f_2(x) = x^4 - x^3 - 4x^2 - 5x - 3$. They have the same degree, so we can take either one of them as the divisor; let's take $f_2(x)$. Divide f_2 into f_1 to get a quotient of 1 and remainder of $f_3(x) = 3x^3 + 4x^2 + 4x + 1$. Then divide f_3 into f_2 to get a quotient and a remainder f_4 , and continue until the remainder is 0 (which occurs on the next iteration).

$$\begin{array}{ll} f_1(x) = x^4 + 2x^3 - x - 2 & f_1(x) = 1 \cdot f_2(x) + f_3(x) \\ f_2(x) = x^4 - x^3 - 4x^2 - 5x - 3 & f_2(x) = \left(\frac{1}{3}x - \frac{7}{9}\right)f_3(x) + f_4(x) \\ f_3(x) = 3x^3 + 4x^2 + 4x + 1 & f_3(x) = \left(\frac{27}{20}x - \frac{9}{20}\right)f_4(x) \\ f_4(x) = -\frac{20}{9}x^2 - \frac{20}{9}x - \frac{20}{9} & \end{array}$$

Thus, a greatest common divisor is $f_4(x)$, which differs by a unit factor from the simpler greatest common divisor $x^2 + x + 1$. We can read the equations on the right in reverse to get f_4 as a linear combination of f_1 and f_2 .

$$\begin{aligned} f_4(x) &= f_2(x) - \left(\frac{1}{3}x - \frac{7}{9}\right)f_3(x) \\ &= f_2(x) - \left(\frac{1}{3}x - \frac{7}{9}\right)(f_1(x) - f_2(x)) \\ &= \left(\frac{1}{3}x + \frac{2}{9}\right)f_2(x) - \left(\frac{1}{3}x - \frac{7}{9}\right)f_1(x) \end{aligned}$$

3.9 Real and complex polynomial rings $\mathbf{R}[x]$ and $\mathbf{C}[x]$

We know a fair amount about $F[x]$, the ring of polynomials over a field F . It has a division algorithm, so it's a Euclidean domain where the Euclidean valuation is the degree of a polynomial, so it has division and Euclidean algorithms. Since it's Euclidean, it's also a principal ideal domain, and that means irreducible elements are prime. And since it's a PID, it's also a unique factorization domain, that is, every polynomial uniquely factors as a product of irreducible polynomials.

Rather than calling irreducible polynomials prime polynomials, we'll use the term "irreducible polynomial". That's the common practice.

The nonzero prime ideals of $F[x]$ are just the principal ideals (f) generated by irreducible polynomials $f \in F[x]$, and, furthermore, they're maximal ideals, so $F[x]/(f)$ is a field. We've seen examples of this, for instance, $\mathbf{R}[x]/(x^2 + 1) \cong \mathbf{R}[i] = \mathbf{C}$, $\mathbf{Q}[x]/(x^2 - 2) \cong \mathbf{Q}(\sqrt{2})$, and $\mathbf{Z}_3[x]/(x^2 + 1) \cong \mathbf{Z}_3(i)$.

So, irreducible polynomials in $F[x]$ give field extensions of F .

The main question for $F[x]$ is: what are the irreducible polynomials?

We'll start with $\mathbf{C}[x]$ and $\mathbf{R}[x]$ followed by $\mathbf{Q}[x]$ and $\mathbf{Z}[x]$.

3.9.1 $\mathbf{C}[x]$ and the Fundamental Theorem of Algebra

In the 16th century Cardano (1501–1576) and Tartaglia (1500–1557) and others found formulas for roots of cubic and quartic equations in terms of square roots and cube roots. At the time, only positive numbers were completely legitimate, negative numbers were still somewhat mysterious, and the first inkling of a complex number appeared. Incidentally, at this time symbolic algebra was still being developed, and they wrote their equations in words instead of symbols!

Here's an illustration of how complex numbers arose. One of Cardano's cubic formulas gives the solution to the equation $x^3 = cx + d$ as

$$x = \sqrt[3]{d/2 + \sqrt{e}} + \sqrt[3]{d/2 - \sqrt{e}}$$

where $e = (d/2)^2 - (c/3)^3$. Bombelli used this to solve the equation $x^3 = 15x + 4$, which was known to have 4 as a solution, to get the solution

$$x = \sqrt[3]{2 + \sqrt{-121}} + \sqrt[3]{2 - \sqrt{-121}}.$$

Now, $\sqrt{-121}$ is not a real number; it's neither positive, negative, nor zero. Bombelli continued to work with this expression until he found equations that lead him to the solution 4. Assuming that the usual operations of arithmetic held for these "numbers," he determined that

$$\sqrt[3]{2 + \sqrt{-121}} = 2 + \sqrt{-1} \quad \text{and} \quad \sqrt[3]{2 - \sqrt{-121}} = 2 - \sqrt{-1}$$

and, therefore, the solution $x = 4$.

Cardano had noted that the sum of the three solutions of a cubic equation $x^3 + bx^2 + cx + d = 0$ is $-b$, the negation of the coefficient of x^2 . By the 17th century the theory of equations had developed so far as to allow Girard (1595–1632) to state a principle of algebra, what we call now "the fundamental theorem of algebra."

His formulation, which he didn't prove, also gives a general relation between the n solutions to an n^{th} degree equation and its n coefficients.

For a generic equation

$$x^n + a_{n-1}x^{n-1} + \cdots + a_1x + a_0 = 0$$

Girard recognized that there could be n solutions, if you allow all roots and count roots with multiplicity. So, for example, the equation $x^2 + 1 = 0$ has the two solutions $\sqrt{-1}$ and $-\sqrt{-1}$, and the equation $x^2 - 2x + 1 = 0$ has the two solutions 1 and 1. Girard wasn't particularly clear what form his solutions were to have, just that there were n of them: x_1, x_2, \dots, x_n .

Girard gave the relation between the n roots x_1, x_2, \dots, x_n and the n coefficients a_1, \dots, a_n that extended Cardano's remark. First, the sum of the roots $x_1 + x_2 + \cdots + x_n$ is $-a_1$ (Cardano's remark). Next, the sum of all products of pairs of solutions is a_2 . Next, the sum of all products of triples of solutions is $-a_3$. And so on until the product of all n solutions is either a_n (when n is even) or $-a_n$ (when n is odd). He figured this out by using a version of one of the properties of polynomials mentioned above, namely, if a_1, a_2, \dots, a_n are roots of a monic polynomial $f(x)$ of degree n , then

$$f(x) = (x - a_1)(x - a_2) \cdots (x - a_n).$$

If you expand the right side of the equation, you'll derive his result.

Here's an example. The 4th degree equation

$$x^4 - 6x^3 + 3x^2 + 26x - 24 = 0$$

has the four solutions $-2, 1, 3,$ and 4 . The sum of the solutions equals 6, that is $-2 + 1 + 3 + 4 = 6$. The sum of all products of pairs (six of them) is

$$(-2)(1) + (-2)(3) + (-2)(4) + (1)(3) + (1)(4) + (3)(4)$$

which is 3. The sum of all products of triples (four of them) is

$$(-2)(1)(3) + (-2)(1)(4) + (-2)(3)(4) + (1)(3)(4)$$

which is -26 . And the product of all four solutions is -24 .

Over the remainder of the 17th century, negative numbers rose in status to be full-fledged numbers. But complex numbers remained suspect through much of the 18th century. They

weren't considered to be real numbers, but they were useful in the theory of equations and becoming more and more useful in analysis. It wasn't even clear what form the solutions to equations might take. Certainly "numbers" of the form $a + b\sqrt{-1}$ were sufficient to solve quadratic equations, even cubic and quartic equations.

Euler did a pretty good job of studying complex numbers. For instance, he studied the unit circle assigning the value $\cos \theta + i \sin \theta$ to the point on the unit circle at an angle θ clockwise from the positive real axis. He measured the angle by the length of the arc cut of the unit circle cut off by the angle. We call that measurement radians now, but the word "radian" wasn't coined until later.

In his study of this circle he developed what we call Euler's identity

$$e^{i\theta} = \cos \theta + i \sin \theta.$$

This was an especially useful observation in the solution of differential equations. Because of this and other uses of i , it became quite acceptable for use in mathematics. By the end of the 18th century numbers of the form $x + iy$ were in fairly common use by research mathematicians, and it became common to represent them as points in the plane.

Yet maybe some other form of "number" was needed for higher-degree equations. The part of the Fundamental Theorem of Algebra which stated there actually are n solutions of an n^{th} degree equation was yet to be proved, pending, of course, some description of the possible forms that the solutions might take.

Still, at nearly the end of the 18th century, it wasn't yet certain what form all the solutions of a polynomial equation might take. Leibniz, for example, stated in 1702 that $x^4 - a^4$ didn't have roots of the form $x + y\sqrt{-1}$, but Euler showed it did in 1742. D'Alembert, Euler, de Foncenex, Lagrange, and Laplace developed partial proofs. Finally, in 1799, Gauss (1777–1855) published his first proof of the Fundamental Theorem of Algebra.

We won't look at his or any other proof of the theorem. That's usually proved in a course in complex analysis. We will, however, use the theorem.

Definition 3.53. A field F is *algebraically closed* if every polynomial $f(x) \in F[x]$ factors as a product of linear factors. Equivalently, a polynomial $f(x)$ of degree n has n roots in F counting multiplicities.

A weaker definition could be made, and that's that every polynomial of degree at least 1 has at least one root in F . By induction, the remaining roots can be shown to exist.

Thus, the Fundamental Theorem of Algebra is a statement that \mathbf{C} is an algebraically closed field. Therefore, the algebra of $\mathbf{C}[x]$ is particularly simple. The irreducible polynomials are the linear polynomials.

3.9.2 The polynomial ring $\mathbf{R}[x]$

Let's turn our attention now to polynomials with real coefficients. Much of what we can say about $\mathbf{R}[x]$ comes from the relation of \mathbf{R} as a subfield \mathbf{C} , and consequently= from the relation of $\mathbf{R}[x]$ as a subring of $\mathbf{C}[x]$. That is to say, we can interpret a polynomial $f(x)$ with real coefficients as a polynomial with complex coefficients.

Theorem 3.54. If a polynomial $f(x)$ with real coefficients has a complex root z , then its complex conjugate \bar{z} is also a root.

Proof. Let $f(x) = a_n x^n + \cdots + a_1 x + a_0$ where each $a_i \in \mathbf{R}$. If z is a root of f , then $f(z) = a_n z^n + \cdots + a_1 z + a_0 = 0$. Take the complex conjugate of the equation, and note that $\bar{a}_i = a_i$. Then $f(\bar{z}) = a_n \bar{z}^n + \cdots + a_1 \bar{z} + a_0 = 0$. Thus, \bar{z} is also a root. Q.E.D.

This theorem tells us for a polynomial $f(x)$ with real coefficients, its roots either come in k pairs of a complex number or singly as real numbers. We can name the $2k$ complex roots as

$$z_1, \bar{z}_1, z_2, \bar{z}_2, \dots, z_k, \bar{z}_k.$$

Writing $z_1 = x_1 + y_1 i, \dots, z_k = x_k + y_k i$, the complex roots are

$$x_1 + y_1 i, x_1 - y_1 i, x_2 + y_2 i, x_2 - y_2 i, \dots, x_k + y_k i, x_k - y_k i$$

and the $n - 2k$ real roots as

$$r_{2k+1}, \dots, r_n.$$

Using the fact that \mathbf{C} is algebraically closed, we can write $f(x)$ as

$$\begin{aligned} f(x) &= a_n(x - z_1)(x - \bar{z}_1) \cdots (x - z_k)(x - \bar{z}_k)(x - r_{2k+1}) \cdots (x - r_n) \\ &= a_n(x^2 - 2x_1 x + x_1^2 + y_1^2) \cdots (x^2 - 2x_k x + x_k^2 + y_k^2)(x - r_{2k+1}) \cdots (x - r_n) \end{aligned}$$

This last expression has factored $f(x)$ as a product of irreducible quadratic and linear polynomials with real coefficients.

Theorem 3.55. The irreducible polynomials in $\mathbf{R}[x]$ are the linear polynomials and the quadratic polynomials with negative discriminant.

Proof. The remarks above show that only linear and quadratic polynomials can be irreducible. Linear polynomials are always irreducible. A quadratic polynomial will have no real roots when its discriminant is negative. Q.E.D.

3.10 Rational and integer polynomial rings

We've studied the irreducible polynomials in $\mathbf{C}[x]$ and $\mathbf{R}[x]$ with the help of the Fundamental Theorem of Algebra and found them to be easily classified. The irreducible polynomials in $\mathbf{C}[x]$ are the linear polynomials, and irreducible polynomials in $\mathbf{R}[x]$ are the linear polynomials and quadratic polynomials with negative discriminant.

Determining which polynomials in $\mathbf{Q}[x]$ are irreducible is much harder. Of course, all the linear ones are, and we'll be able to tell which quadratic and cubic ones are irreducible fairly easily. After that it becomes difficult.

3.10.1 Roots of polynomials

The quadratic case. Consider a quadratic polynomial $f(x) = ax^2 + bx + c$ with coefficients in \mathbf{Q} .

Its roots are given by the quadratic formula $\frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$ which can be shown by the process known as completing the square. The *discriminant* of a quadratic polynomial is

$\Delta = b^2 - 4ac$. When Δ is positive, there are two real roots; when 0, there is one double root; and when negative, the roots are a pair of complex conjugate numbers.

When Δ is a perfect rational square, that is, the square of a rational number, then $f(x)$ factors, that is, it's reducible. Otherwise, it's irreducible.

Thus, $f(x)$ is irreducible if and only if the discriminant is not a perfect square.

The cubic case. It is more difficult to determine when a cubic polynomial $f(x) = ax^3 + bx^2 + cx + d$ with rational coefficients is irreducible, but not too difficult. Note that if $f(x)$ factors, then one of the factors has to be linear, so the question of reducibility reduces to the existence of a rational root of $f(x)$.

Various solutions of a cubic equation $ax^3 + bx^2 + cx + d = 0$ have been developed. Here's one. First, we may assume that f is monic by dividing by the leading coefficient. Our equation now has the form $x^3 + bx^2 + cx + d = 0$. Second, we can eliminate the quadratic term by replacing x by $y - \frac{1}{3}b$. The new polynomial in y will have different roots, but they're only translations by $\frac{1}{3}b$. We now have the cubic equation

$$y^3 + (c - \frac{1}{3}b^2)y + (\frac{2}{27}b^3 - \frac{1}{3}bc + d) = 0$$

which we'll write as

$$y^3 + py + q = 0.$$

By the way, this substitution which results in a polynomial whose term after the leading term is 0 has a name. It is called a Tschirnhaus substitution. The roots of the new polynomial will sum to 0.

We'll follow Viète's method and perform another substitution. Replace y by $z - \frac{p}{3z}$. After simplifying and clearing the denominators we'll have the equation

$$z^6 + qz^3 - \frac{p^3}{27z} = 0$$

which is a quadratic equation in z^3 . Its complex solutions are

$$z^3 = \frac{-q \pm \sqrt{q^2 + 4p^3/27}}{2} = -\frac{q}{2} \pm \sqrt{\left(\frac{q}{2}\right)^2 + \left(\frac{p}{3}\right)^3}.$$

Taking complex cube roots to get three values for z , then using $y = z - \frac{p}{3z}$ to determine y and $x = y - \frac{1}{3}b$ to determine x , we have the all three complex solutions to the original equation. At least one of these three complex solutions is real, and perhaps all three.

We have a way of determining whether a cubic polynomial is reducible. First $q^2 + 4p^3/27$ needs to be a perfect rational square r^2 , then one of $-q + r$ and $-q - r$ needs to be a perfect rational cube.

There is another way to determine if there is a rational root.

Rational roots of a polynomial. If we're looking for the roots of a polynomial with rational coefficients, we can simplify the job a little bit by clearing the denominators so that all the coefficients are integers. The following theorem helps in finding roots.

Theorem 3.56 (Rational root theorem). Let $f(x) = a_n x^n + \cdots + a_1 x + a_0$ be a polynomial with integral coefficients. If r/s is a rational root of f with r/s in lowest terms, then r divides the constant a_0 and s divides the leading coefficient a_n .

Proof. Since r/s is a root, therefore

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

and so, clearing the denominators, we have

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

We can rewrite this equation as

$$(a_n r^{n-1} + a_{n-1} r^{n-2} s + \cdots + a_1 s^{n-1}) r = -a_0 s^n.$$

Now, since r divides $-a_0 s^n$, and r is relatively prime to s , and hence to s^n , therefore r divides a_0 . In like manner, you can show s divides a_n . Q.E.D.

For example, to find the rational roots r/s of $f(x) = 27x^4 + 30x^3 + 26x^2 - x - 4$, r will have to divide 4, so the possibilities for r are $\pm 1, \pm 2, \pm 4$, and s will have to divide 27, so the possibilities for s are 1, 3, 9, 27 (since we may assume s is positive). That gives 24 rational numbers to check, and among them will be found the two rational roots $\frac{1}{3}$ and $-\frac{4}{9}$. After one, $\frac{r}{s}$, is found $f(x)$ can be divided by $x - \frac{r}{s}$ to lower the degree of the polynomial to find the rest of the roots.

If a polynomial does have a rational root, then it's clearly reducible since that rational root determines a linear factor of the polynomial. That gives us another way to determine if a cubic polynomial is reducible.

For polynomials of degree 4 or higher, knowing that there are no rational roots is insufficient to conclude the polynomial is irreducible. It still may factor as quadratic and higher degree terms. For example, $x^4 + x^2 + 1$ has no rational roots, but it factors as $(x^2 + x + 1)(x^2 - x + 1)$, so it is reducible.

3.10.2 Gauss's lemma and Eisenstein's criterion

Further study of $\mathbf{Q}[x]$ will require looking at $\mathbf{Z}[x]$. In other words, in order to study polynomials with rational coefficients, we'll have to look at polynomials with integral coefficients. We can take a polynomial with rational coefficients and multiply it by the least common multiple of the denominators of its coefficients to get another polynomial with the same roots but with integral coefficients. We can also divide by the greatest common divisor of the resulting coefficients to get yet another polynomial with the same roots, with integral coefficients, and the greatest common divisor of all its coefficients is 1. Such a polynomial is called *primitive*.

After that, we'll be able to prove Gauss's lemma which says that a primitive polynomial $f(x) \in \mathbf{Z}[x]$ is reducible in $\mathbf{Q}[x]$ if and only if it's reducible in $\mathbf{Z}[x]$.

We can make more use of these results if, instead of considering just the case of the domain \mathbf{Z} and its field of fractions \mathbf{Q} , we generalize to any unique factorization domain D and its field of fractions F . So, for the following discussion, fix a UFD D , and let F denote its field

of fractions. Though, keep in mind the basic case when $D = \mathbf{Z}$, $F = \mathbf{Q}$, $D/(p) = \mathbf{Z}_p$, and $D/(p)[x] = \mathbf{Z}_p[x]$ to get a better idea of what's going on.

When we have a prime p in D , the projection $\gamma : D \rightarrow D/(p)$ induces a ring epimorphism $D[x] \rightarrow D/(p)[x]$ between polynomial rings where the coefficients of f are reduced modulo p giving a polynomial in $D/(p)[x]$. We'll denote the resulting polynomial in $D/(p)[x]$ by f_p .

Definition 3.57. The *content* of a polynomial in $D[x]$ is the greatest common divisor of all of its coefficients. If the content is 1, the polynomial is called *primitive*.

For example, if $f(x) = 3x^2 - 9x + 6$ then the content of f is 3. Also, the content of a monic polynomial is 1, so all monic polynomials are primitive.

The content of a polynomial is only defined up to a unit.

Evidently, every polynomial in $D[x]$ equals a constant times a primitive polynomial, the constant being its content.

Lemma 3.58 (Gauss). The product of two primitive polynomials in $D[x]$ is primitive, and the content of the product of any two polynomials in $D[x]$ is the product of their contents (up to a unit).

Proof. In order to show the first statement, we'll show if the product is not primitive, then one of the two polynomials is not primitive.

Let f and g be primitive polynomials and suppose that their product fg is not primitive. Then some prime p of D divides the content of fg , so p divides every coefficient of fg . Therefore, in $D/(p)[x]$, $(fg)_p = 0$, so $f_p g_p = 0$. But $D/(p)[x]$ is an integral domain (in fact, a UFD), so either $f_p = 0$ or $g_p = 0$. Therefore, p either divides all the coefficients of f or all the coefficients of g , hence one or the other is not primitive.

The second statement follows from the first just by using the fact that a polynomial equals its content times a primitive polynomial. Q.E.D.

Theorem 3.59 (Gauss's lemma). If a primitive polynomial in $D[x]$ can be factored as the product of two polynomials in $F[x]$, then it can be factored as the product of two polynomials in $D[x]$ of the same degrees.

Proof. Given $f \in D[x]$ as a product gh with $g, h \in F[x]$. We can write $gh = \frac{p}{q} uv$ where u and v are primitive polynomials in $D[x]$, and p and q are relatively prime integers. Then $qf = puv$. Since f is primitive, the content of qf equals the content of q . Since u and v are primitive, so is uv , and therefore the content of puv equals the content of p . Thus $p = q$, and they're both 1, and so $f = uv$. Note that the degrees of u and v are the same as the degrees of g and h , respectively. Q.E.D.

The following corollary is sometimes called Gauss's lemma. It follows directly from the above since monic polynomials are primitive.

Corollary 3.60. A monic polynomial in $D[x]$ is reducible over $F[x]$ if and only if it's reducible over $D[x]$.

There are irreducibility tests for polynomials with integer coefficients, so by this corollary, we'll be able to test irreducibility for polynomials with rational coefficients.

One test for irreducibility of polynomials with integer coefficients is to move to a quotient ring \mathbf{Z}_p . That also generalizes to any UFD D . If you can factor it in D , you can factor it in a quotient ring, at least if the leading term doesn't disappear in the quotient.

Theorem 3.61 (Modulo p irreducibility test.). Let p be a prime integer, and let f be a polynomial whose leading coefficient is not divisible by p . If f is reducible in $F[x]$, then f_p is reducible in $D/(p)[x]$. If f_p is irreducible in $D/(p)[x]$, then f is irreducible in $F[x]$.

Proof. Suppose f is reducible in $F[x]$. Then there exist $g, h \in D[x]$ such that $f = gh$ where the degrees of g and h are at least 1. Since $f = gh$, therefore, $f_p = g_p h_p$. Since p does not divide the leading coefficient of f , neither does it divide the leading coefficients of g or h . Therefore $\deg g_p = \deg g \geq 1$ and $\deg h_p = \deg h \geq 1$. Thus, f_p is reducible.

The last statement of the theorem is the contrapositive of the first statement. Q.E.D.

Example 3.62. Consider any cubic polynomial f in $\mathbf{Q}[x]$ with an odd leading coefficient, an odd constant, and one of the other two coefficients odd, for instance, $f(x) = 77x^3 + 15x^2 + 8x + 105$. By Gauss's lemma, it's reducible in $\mathbf{Q}[x]$ if and only if it's reducible in $\mathbf{Z}[x]$. To determine that, use the modulo 2 irreducible test. For $f(x) = 77x^3 + 15x^2 + 8x + 105$, you'll get $f_2(x) = x^3 + x^2 + 1$. The resulting f_2 will have no roots in \mathbf{Z}_2 since it has three nonzero terms. A cubic polynomial with no roots is irreducible, so f_2 is irreducible in $\mathbf{Z}_2[x]$. Hence, f is irreducible in $\mathbf{Q}[x]$.

The converse of the mod p irreducibility test is not valid. A polynomial can be reducible mod p but irreducible in $\mathbf{Z}[x]$. Take $f(x) = 77x^3 + 15x^2 + 8x + 105$, for example, which we know is irreducible in $\mathbf{Z}[x]$. Modulo $p = 5$, however, it factors into linear factors: $f_5(x) \equiv 2x^3 - 2x = 2(x + 1)x(x - 1)$, so is reducible.

Exercise 44. Show the polynomial $f(x) = x^4 + x^3 + 2x^2 + 2x + 1$ is irreducible in $\mathbf{Q}[x]$. Hint: Consider it modulo 2. First check for roots, then see if it's divisible by a irreducible quadratic. There aren't many irreducible quadratics modulo 2; $x^2 + 1$ isn't since it factors as $(x + 1)^2$ modulo 2. Neither are x^2 or $x^2 + x$ since they're both divisible by x .

Another useful irreducibility test is Eisenstein's criterion.

Theorem 3.63 (Eisenstein's criterion). Let $f \in D[x]$. If a prime p does not divide the leading coefficient of f , but it does divide all the other coefficients, and p^2 does not divide the constant of f , then f is irreducible in $F[x]$.

Proof. Suppose f is reducible. As in the previous theorem, there exist $g, h \in D[x]$ such that $f = gh$ where the degrees of g and h are at least 1. Reduce everything modulo p . Then $a_n x^n = f_p(x) = g_p(x)h_p(x)$ where a_n is the leading coefficient of f . Now $\mathbf{Z}_p[x]$ is a UFD, and since $f_p(x)$ is the unit a_n times the irreducible x raised to the n^{th} power, therefore x divides both $g_p(x)$ and $h_p(x)$. Therefore $g_p(0) = h_p(0) = 0$. That means that p divides the constant terms of both g and h , which implies p^2 divides the constant term of f , contrary to the assumption. Q.E.D.

Example 3.64. Consider the polynomial $f(x) = x^n - a$. As long as a has a prime factor that appears to the first power, then Eisenstein's criterion implies f is irreducible.

Exercise 45. Show that the polynomial $f(x) = x^n + 10x + 15$ is irreducible in $\mathbf{Q}[x]$

3.10.3 Prime cyclotomic polynomials

Cyclotomic polynomials were introduced in definition 1.63.

For a prime p , the p^{th} cyclotomic polynomial is

$$\Phi_p(x) = \frac{x^p - 1}{x - 1} = x^{p-1} + \cdots + x + 1.$$

We'll use Eisenstein's criterion to show Φ_p is irreducible, but not directly. First, we'll use a translation. Let

$$f(x) = \Phi_p(x + 1) = \frac{(x + 1)^p - 1}{x} = x^{p-1} + \binom{p}{p-1}x^{p-2} + \cdots + \binom{p}{2}x + \binom{p}{1}.$$

Then Eisenstein's criterion applies to f . Since f is irreducible, so is Φ_p .

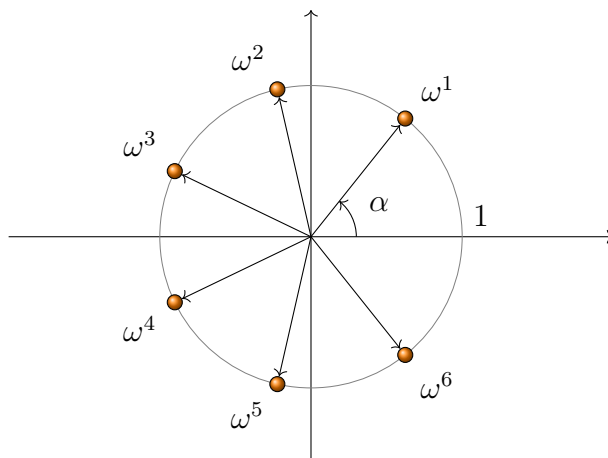


Figure 3.4: Primitive 7th roots of unity

The roots of Φ_p are the $p - 1$ primitive p^{th} roots of unity. In the case that $p = 7$, there are six primitive 7th roots of unity. One, labelled ω in figure 3.4, is located at an angle of $\alpha = 2\pi/7$, and the others are powers of it. There is one more 7th root of unity, namely 1. 1 is not primitive 7th root of unity but instead a primitive first root of unity since it's a root of the polynomial $x - 1$ of degree 1.

3.10.4 Polynomial rings with coefficients in a UFD, and polynomial rings in several variables.

Gauss's lemma has more uses than we've used it for. We can use it to show that if D is a UFD, then so is the polynomial ring $D[x]$. And we can apply that statement to conclude a polynomial ring $D[x, y]$ in two or $D[x_1, \dots, x_n]$ more variables is also a UFD. Although these rings are UFDs, they're not PIDs.

Theorem 3.65. Let D be a unique factorization domain and F its ring of fractions. Then $D[x]$ is also a UFD. The irreducible polynomials in $D[x]$ are either irreducible elements of D or have content 1 and are irreducible polynomials in $F[x]$.

Proof. Let f be a nonzero polynomial in $D[x]$. It is equal to its content times a primitive polynomial. Its content is an element of D , and, since D is a UFD, its content uniquely factors (up to a unit) as a product of irreducible elements of D .

We're reduced the proof to showing that that a primitive polynomial f in $D[x]$ of degree at least 1 uniquely factors as a product of irreducible polynomials.

Since f is a polynomial in $D[x]$, it's also a polynomial in $F[x]$, and we know $F[x]$ is a UFD being a polynomial ring with coefficients in a field F . Thus, f uniquely factors in F :

$$f(x) = f_1(x)f_2(x) \cdots f_k(x)$$

where each $f_i(x)$ is irreducible in $F[x]$. We only need to show that this factorization can be carried out in $D[x]$. Each polynomial $f_i(x)$ is a element a_i of F times a primitive polynomial $f'_i(x)$ in $D[x]$, so

$$f(x) = a_1 \cdots a_k f'_1(x) \cdots f'_k(x).$$

Since $f(x)$ is primitive and the product $f'_1(x) \cdots f'_k(x)$ is also primitive, therefore $a_1 \cdots a_k$ is a unit in D . Thus, $f(x)$ factors in $D[x]$. You can also show that it can factor in only one way in $D[x]$ since it only factors in one way in $F[x]$. Q.E.D.

Corollary 3.66. If D is a UFD, then a polynomial ring in several variables $D[x_1, x_2, \dots, x_r]$ with coefficients in D is also a UFD.

In general, these aren't PIDs. For example, $(2, x)$ is not a principal ideal in $\mathbf{Z}[x]$, and (x, y) is not a principal ideal in $\mathbf{Q}[x, y]$.

Irreducible polynomials and field extensions. We'll see irreducible polynomials in a field $F[x]$ correspond to maximal ideals. Quotient rings by ideals are fields if and only if the ideal is a maximal ideal, as shown in section 3.6.4. Therefore irreducible polynomials correspond to field extensions.

Theorem 3.67. The ideal generated by a polynomial with coefficients in a field is maximal if and only if the polynomial is irreducible over the field.

Proof. Let $f \in F[x]$. Suppose first that (f) is a maximal ideal of $F[x]$. We'll show that f can't be reducible. Suppose that f factors as gh in $F[x]$ where g and h have lower degrees than f , but neither g nor h is a unit. Then $(f) \subseteq (g) \subseteq F[x]$. Since (f) is maximal, therefore (g) is equal to either (f) or $F[x]$. Neither of these can occur, for if $(g) = (f)$, then they have the same degree; but if $(g) = F[x]$, then g is a unit. Therefore f is not a reducible polynomial.

Next suppose that f is an irreducible polynomial in $F[x]$. Let I be any ideal such that $(f) \subseteq I \subset F[x]$. Since $F[x]$ is a Euclidean domain (see section 3.8.4), it is also a principal ideal domain, so $I = (g)$ for some polynomial g . Therefore $f = gh$ for some $h \in F[x]$. But f is irreducible, so either g or h is a unit. But g is not a unit since $(g) = I \neq F[x]$, so h is a unit. Therefore $(f) = (g)$. So any proper ideal I that contains (f) is (f) itself. Thus (f) is a maximal ideal. Q.E.D.

Corollary 3.68. The quotient ring $F[x]/(f)$ of a polynomial ring over a field F by the ideal generated by an irreducible polynomial is a field extension of F .

This corollary follows directly the preceding theorem and the maximal ideal theorem, theorem 3.36, .

Example 3.69. We saw earlier that any cubic polynomial f in $\mathbf{Q}[x]$ with an odd leading coefficient, an odd constant, and one of the other two coefficients odd is irreducible. So, for example, $f(x) = x^3 - x - 1$ is irreducible over \mathbf{Q} . That means $K = \mathbf{Q}[x]/(f)$ is a field. This field K can also be denoted $K = \mathbf{Q}[x]/(x^3 = x + 1)$ since in the quotient, $x^3 - x - 1$ is 0. Rather than using the symbol x in the quotient, it would be better to have a different symbol so that x can still be used as our variable. Let's use w .

Then every element in it is of the form $aw^2 + bw + c$. Addition is done as usual for polynomials, and multiplication is as usual except whenever w^3 appears, it is replaced by $w + 1$. For example, the product $(w^2 + 3w - 3)(2w^2 - 5) = 2w^3 + w^2 - 9w + 15 \equiv w^2 - 7w + 17$.

As this is a field, there are reciprocals of nonzero elements and division by nonzero elements. Finding reciprocals is not so easy. For example, to find the reciprocal of w , we need an element $aw^2 + bw + c$ such that its product with w equals 1. Now, $(aw^2 + bw + c)w = aw^3 + bw^2 + cw = bw^2 + (a + c)w + a$, so for that to equal 1, $a = 1$, $b = 0$, and $c = -1$. So the reciprocal is $w^{-1} = w^2 - 1$.

Exercise 46. Find an irreducible cubic polynomial in $\mathbf{Z}_2[x]$ to construct a field with eight elements. Write down a multiplication table for that field. You can leave out 0 and 1 from the table since it's obvious how they multiply, but have six rows and columns labeled $a = x$, $b = x + 1$, $c = x^2$, $d = x^2 + 1$, $e = x^2 + x$, and $f = x^2 + x + 1$.

3.11 Number fields and their rings of integers

In section 2.3 a *number field* K was defined a finite field extension of the rational number field \mathbf{Q} . Some examples were $\mathbf{Q}[\sqrt{2}]$ and $\mathbf{Q}[\sqrt{-1}] = \mathbf{Q}[i]$.

The Gaussian integers $\mathbf{Z}[i]$ is an example of what is called a ring of integers. We'll see in this section what a ring of integers is and study some of their properties.

A number field was defined to be an algebraic extension of \mathbf{Q} , and an algebraic integer was defined in section 2.3.1 to be a root of a monic polynomial with coefficients in the integers.

Our first goal is to show that the set of all algebraic integers in number field K , that set being denoted \mathcal{O}_K is a subring of K . That will take a few steps.

Define the minimal polynomial of an algebraic integer a is being that monic polynomial f in \mathbf{Q} of minimal degree such that $f(a) = 0$. Note that we're not requiring f to have coefficients in \mathbf{Z} , but we'll prove that below.

Lemma 3.70. The minimal polynomial of an algebraic integer divides every polynomial in $\mathbf{Q}[x]$ of which it is a root.

Proof. Let a be an algebraic integer with minimal polynomial f . By the division algorithm, there are polynomials q and r in $\mathbf{Q}[x]$ such that $g = qf + r$, where either $r = 0$ or $\deg r < \deg f$. Then $r(a) = g(a) - q(a)f(a) = 0$, so a is a root of r . Since f is the polynomial of least positive degree with root a , so $r = 0$. Q.E.D.

Lemma 3.71. The minimal polynomial of an algebraic integer has coefficients in \mathbf{Z} .

Proof. Let f be the minimal polynomial of a , and let g be a monic polynomial in $\mathbf{Z}[x]$ such that $g(a) = 0$. By the previous lemma, $g = fh$ for some $h \in \mathbf{Q}[x]$.

Suppose that $f \notin \mathbf{Z}[x]$, then some prime number p divides the denominator of some coefficient of f . Let p^i be the largest power of p dividing that denominator, so $i \geq 1$. Let p^j be the largest power of p that divides some denominator of a coefficient of h , with $j \geq 0$. Then $p^{i+j}g = (p^i f)(p^j h)$. Now take that equation modulo p . Modulo p , the left side is 0, but neither polynomial on the right side is 0, a contradiction since $\mathbf{Z}_p[x]$ is an integral domain. Q.E.D.

Theorem 3.72. The set of all algebraic integers, \mathcal{O}_K , in a number field K is a subring of that field.

This ring \mathcal{O}_K is called *the ring of integers* in the number field K .

Exercise 47. Prove that the Gaussian integers $\mathbf{Z}[i]$ is the ring of integers in the number field $\mathbf{Q}[i]$. Since i is the root of the monic polynomial $x^2 + 1$, all that's needed to prove is that there are no other integral elements in $\mathbf{Q}[i]$ other than those in $\mathbf{Z}[i]$.

Chapter 4

Groups

Recall that a group is a set equipped with one binary operation that is associative, has an identity element, and has inverse elements. If that binary operation is commutative, then the group is called an Abelian group.

4.1 Groups and subgroups

4.1.1 Definition and basic properties of groups

We'll look at basic properties of groups, and since we'll discuss groups in general, we'll use a multiplicative notation even though some of the example groups are Abelian.

Definition 4.1. The axioms for a group are very few. A group G has an underlying set, also denoted G , and a binary operation $G \times G \rightarrow G$ that satisfies three properties.

1. Associativity. $(xy)z = x(yz)$.
2. Identity. There is an element 1 such that $1x = x = x1$.
3. Inverses. For each element x there is an element x^{-1} such that $xx^{-1} = x^{-1}x = 1$.

Theorem 4.2. From these few axioms several properties of groups immediately follow.

1. Uniqueness of the identity. There is only one element e such that $ex = x = xe$, and it is $e = 1$.

Outline of proof. The definition says that there is at least one such element. To show that it's the only one, suppose e also has the property of an identity and prove $e = 1$.

2. Uniqueness of inverses. For each element x there is only one element y such that $xy = yx = 1$.

Outline of proof. The definition says that there is at least one such element. To show that it's the only one, suppose that y also has the property of an inverse of x and prove $y = x^{-1}$.

3. Inverse of an inverse. $(x^{-1})^{-1} = x$.

Outline of proof. Show that x has the property of an inverse of x^{-1} and use the previous result.

4. Inverse of a product. $(xy)^{-1} = y^{-1}x^{-1}$.
Outline of proof. Show that $y^{-1}x^{-1}$ has the property of an inverse of xy .
5. Cancellation. If $xy = xz$, then $y = z$, and if $xz = yz$, then $x = y$.
6. Solutions to equations. Given elements a and b there are unique solutions to each of the equations $ax = b$ and $ya = b$, namely, $x = a^{-1}b$ and $y = ba^{-1}$.
7. Generalized associativity. The value of a product $x_1x_2 \cdots x_n$ is not affected by the placement of parentheses.
Outline of proof. The associativity in the definition of groups is for $n = 3$. Induction is needed for $n > 3$.
8. Powers of an element. You can define x^n for nonnegative values of n inductively. For the base case, define $x^0 = 1$, and for the inductive step, define $x^{n+1} = xx^n$. For negative values of n , define $x^n = (x^{-n})^{-1}$.
9. Properties of powers. Using the definition above, you can prove using induction the following properties of powers where m and n are any integers: $x^m x^n = x^{m+n}$, $(x^m)^n = x^{mn}$.

Note that $(xy)^n$ does not equal $x^n y^n$ in general, although it does for Abelian groups.

4.1.2 Subgroups

A subgroup H of G is a group whose underlying set is a subset of the underlying set of G and has the same binary operation, that is, for $x, y \in H$, $x \cdot_H y = x \cdot_G y$ where \cdot_H denotes is the multiplication in H while \cdot_G denotes is the multiplication in G . Since they are the same, we won't have to subscript the multiplication operation.

An alternate description of a subgroup H is that it is a subset of G that is closed under multiplication, has 1, and is closed under inverses.

Of course, G is a subgroup of itself. All other subgroups of G , that is, those subgroups that don't have every element of G in them, are called *proper subgroups*.

Also, $\{1\}$ is a subgroup of G , usually simply denoted 1. It's called the *trivial subgroup* of G .

Example 4.3. Consider the cyclic group of six elements $G = \{1, a, a^2, a^3, a^4, a^5\}$ where $a^6 = 1$. Besides the trivial subgroup 1 and the entire subgroup G , there are two other subgroups of G . One is the 3-element subgroup $\{1, a^2, a^4\}$ and the other is the 2-element subgroup $\{1, a^3\}$.

The intersection $H \cap K$ of two subgroups H and K is also a subgroup, as you can easily show. Indeed, the intersection of any number of subgroups is a subgroup.

The union of two subgroups is never a subgroup unless one of the two subgroups is contained in the other.

Exercise 48. About intersections and unions of subgroups.

- (a). Show that the intersection of two subgroups is also a subgroup.
- (b). Give a counterexample where the union of two subgroups is not a subgroup.

Example 4.4 (Subgroups of \mathbf{Z}). Consider the group \mathbf{Z} under addition. A subgroup of \mathbf{Z} has to be closed under addition, include 0, and be closed under negation. Besides 0 and \mathbf{Z} itself, what are the subgroups of \mathbf{Z} ? If the subgroup is nontrivial, then it has a smallest positive element, n . But if n lies in a subgroup, then all multiples, both positive and negative, of n also must be in the subgroup. Thus, $n\mathbf{Z}$ is that subgroup of \mathbf{Z} .

Useful subgroups of a group. There are a number of other subgroups of a group that are important in studying nonabelian groups such as the center of a group and the centralizer of an element of a group.

Definition 4.5. Center and centralizer.

The *center* of a group G is $Z(G) = \{x \in G \mid ax = xa \text{ for all } a \in G\}$,

For $a \in G$, the *centralizer* of a is $Z_a(G) = \{x \in G \mid ax = xa\}$.

Exercise 49. Show the following properties about centers and centralizers.

- Prove that $Z(G)$ is a subgroup of G .
- Prove that the center of G is the intersection of all the centralizer subgroups of G .
- Prove that $Z_a(G)$ is a subgroup of G .

Definition 4.6 (Commutator subgroup). The commutator of two elements x and y in a group G is the element $x^{-1}y^{-1}xy$. It is denoted $[x, y]$.

The subgroup of G generated by all the commutators of its elements is called the *commutator subgroup* of G , denoted G' .

Note that for an Abelian group, all the commutators are 1, and the the commutator subgroup is trivial.

If S is a subset of G , then there is a smallest subgroup $\langle S \rangle$ of G containing S . It can be described as the intersection of all subgroups H containing S ,

$$\langle S \rangle = \bigcap_{S \subseteq H} H.$$

Alternatively, it can be described as the subset of G of all products of powers of elements of S ,

$$\langle S \rangle = \{x_1^{e_1} x_2^{e_2} \cdots x_n^{e_n} \mid n \geq 0, \text{ each } x_i \in S, \text{ and each } e_i \in \mathbf{Z}\}.$$

4.1.3 Cyclic groups and subgroups

If a is an element of a group G , then the subset of G generated by a

$$\langle a \rangle = \{a^n \mid n \in \mathbf{Z}\}$$

is a subgroup of G . This subgroup generated by a is called a *cyclic subgroup* of G . If G itself is generated by some element a , then G is called a *cyclic group*.

Definition 4.7 (Order and involution). The *order* of a group G is the number of elements in it, that is, the cardinality of its underlying set. It's usually denoted $|G|$.

The *order* of an element a in a group is the smallest positive integer n such that $a^n = 1$. It's denoted $\text{ord } a$. If every positive power $a^n \neq 1$, then the order of n is defined to be ∞ . So, for example, the order of 1 is 1 since $1^1 = 1$.

An *involution* a is an element of a group which is its own inverse, $a^{-1} = a$. Clearly, the order of an involution a is 2 unless $a = 1$, in which case the order of a is 1.

Exercise 50. Prove that the order of a is also equal to the order of the cyclic group $\langle a \rangle$ generated by a . That is, $\text{ord } a = |\langle a \rangle|$.

An abstract cyclic group of order n is often denoted $C_n = \{1, a, a^2, \dots, a^{n-1}\}$ when the operation is written multiplicatively. It is isomorphic to the underlying additive group of the ring \mathbf{Z}_n where an isomorphism is $f : \mathbf{Z}_n \rightarrow C_n$ is defined by $f(k) = a^k$.

Exercise 51. Prove that any subgroup of a cyclic group is itself cyclic.

Exercise 52. Let G be a cyclic group of order n and a an element of G . Prove that a generates G , that is, $\langle a \rangle = G$, if and only if $\text{ord } a = n$.

Cyclic groups are all Abelian, since $a^n a^m = a^{m+n} = a^m a^n$. The integers \mathbf{Z} under addition is an infinite cyclic group, while \mathbf{Z}_n , the integers modulo n , is a finite cyclic group of order n .

Exercise 53. Prove that every cyclic group is isomorphic either to \mathbf{Z} or to \mathbf{Z}_n for some n .

Exercise 54. Prove that if k is relatively prime to n , then k generates \mathbf{Z}_n .

4.1.4 Products of groups

Just as products of rings are defined coordinatewise, so are products of groups. Using multiplicative notation, if G and H are two groups then $G \times H$ is a group where the product $(x_1, y_1)(x_2, y_2)$ is defined by $(x_1 x_2, y_1 y_2)$. The identity element in $G \times H$ is $(1, 1)$, and the inverse $(x, y)^{-1}$ is (x^{-1}, y^{-1}) . The projections $\pi_1 : G \times H \rightarrow G$ and $\pi_2 : G \times H \rightarrow H$ are group epimorphisms where $\pi_1(x, y) = x$ and $\pi_2(x, y) = y$.

Also, $\iota_1 : G \rightarrow G \times H$ and $\iota_2 : H \rightarrow G \times H$ are group monomorphisms where $\iota_1(x) = (x, 1)$ and $\iota_2(y) = (1, y)$. Thus, we can interpret G and H as subgroups of $G \times H$.

Note that G and H are both Abelian groups if and only if $G \times H$ is an Abelian group. The product of two Abelian groups is also called their *direct sum*, denoted $G \oplus H$.

The underlying additive group of a ring is an Abelian group, and some of the results we have for rings give us theorems for Abelian groups. In particular, the Chinese remainder theorem for cyclic rings \mathbf{Z}_n gives us a theorem for cyclic groups C_n .

Theorem 4.8 (Chinese remainder theorem for groups). Suppose that $n = km$ where k and m are relatively prime. Then the cyclic group C_n is isomorphic to $C_k \times C_m$. More generally, if n is the product $k_1 \cdots k_r$ where the factors are pairwise relatively prime, then

$$C_n \cong C_{k_1} \times \cdots \times C_{k_r} = \prod_{i=1}^r C_{k_i}.$$

In particular, if the prime factorization of n is $n = p_1^{e_1} \cdots p_r^{e_r}$. Then the cyclic group C_n factors as the product of the cyclic groups $C_{p_i^{e_i}}$, that is,

$$C_n \cong \prod_{i=1}^r C_{p_i^{e_i}}.$$

4.1.5 Cosets and Lagrange's theorem

Cosets are useful in developing the combinatorics of finite groups, that is, for counting subgroups and other things related to a finite group. They come in both left and right forms as you'll see in the definition below, but we'll only use left cosets. Our first combinatorial theorem is called Lagrange's theorem which says that the order of a subgroup divides the order of a group. Since the subgroup $\langle a \rangle$ generated by a single element has an order that divides the order of the group, therefore the order of an element divides the order of the group, too. We'll have our first classification theorem as a corollary, and that is that a group whose order is a prime number is cyclic. Thus, up to isomorphism, there is only one group of that order.

Definition 4.9. Let H be a subgroup of G . A *left coset* of H is a set of the form

$$aH = \{ah \mid h \in H\}$$

while a *right coset* is of the form $Ha = \{ha \mid h \in H\}$.

Theorem 4.10. Several properties of cosets follow from this definition.

1. The coset $1H$ is just the subgroup H itself. In fact, if $h \in H$ then $hH = H$.
2. More generally, $aH = bH$ if and only if $ab^{-1} \in H$. Thus, the same coset can be named in many different ways.
3. Cosets are disjoint. If $aH \neq bH$, then $aH \cap bH = \emptyset$.

Outline of proof. It's probably easier to show the contrapositive: if $aH \cap bH \neq \emptyset$ then $aH = bH$. Suppose an element is in the intersection. Then it can be written as ah or as bh' where both h and h' are elements of H . The rest relies on the previous statement.

4. Cosets of H all have the same cardinality.

Outline of proof. Check that the function $f(ah) = bh$ is a bijection $aH \rightarrow bH$.

5. Thus, the cosets of H partition G into subsets all having the same cardinality.
6. *Lagrange's theorem.* If G is a finite group, and H a subgroup of G , then $|H|$ divides $|G|$. Moreover, $|G|/|H|$ is the number of cosets of H .

Outline of proof. Follows from the preceding statement.

Definition 4.11. The *index* of a subgroup H of a group G is the number of cosets of H . The index is denoted $[G : H]$. By Lagrange's theorem, $[G : H] = |G|/|H|$ when G is a finite group.

Corollary 4.12. If the order of a group is a prime number, then the group is cyclic.

Proof. Let $|G| = p$, a prime. Since p has no divisors except 1 and p , therefore, by Lagrange's theorem, G only has itself and the trivial subgroup as its subgroups. Let $a \neq 1$ be an element of G . It generates a cyclic subgroup $\langle a \rangle$ which isn't trivial, so $\langle a \rangle = G$. Thus G is cyclic. Q.E.D.

Corollary 4.13. If a group is finite, then the order of every element divides the order of the group.

Proof. Let a be an element of a finite group G . Then the order of the subgroup $\langle a \rangle$ divides $|G|$. But $\text{ord } a$ is the order of $\langle a \rangle$. Therefore $\text{ord } a$ divides $|G|$. Q.E.D.

Products of subsets in a group. Occasionally we'll want to look at products HK of subsets H and K , especially when H and K are subgroups of a group G . This product is defined by

$$HK = \{xy \mid x \in H, y \in K\}.$$

Even when H and K are subgroups, it isn't necessary that HK is a subgroup, but there is a simple criterion to test if it is.

Abelian groups are often written additively. In that case, rather than using the notation HK , the notation $H + K$ is preferred: $H + K = \{x + y \mid x \in H, y \in K\}$.

Theorem 4.14. Let H and K be subgroups of G . Then HK is also a subgroup of G if and only if $HK = KH$.

Proof. \implies : Suppose that HK is a subgroup. We'll show that $KH \subseteq HK$. Let $xy \in KH$ with $x \in K$ and $y \in H$. Since $x = 1x \in HK$ and $y = y1 \in HK$, therefore their product xy is also in HK . Thus, $KH \subseteq HK$. Likewise $HK \subseteq KH$. Therefore $HK = KH$.

\impliedby : Suppose that $HK = KH$. To show it's a subgroup, first note $1 \in HK$ since $1 \in H$ and $1 \in K$.

Second, we'll show that HK is closed under multiplication. Let x_1y_1 and x_2y_2 be elements of HK with $x_1, x_2 \in H$ and $y_1, y_2 \in K$. Then $y_1x_2 \in KH = HK$, so $y_1x_2 = x_3y_3$ where $x_3 \in H$ and $y_3 \in K$. Therefore, $(x_1y_1)(x_2y_2) = (x_1x_3)(y_3y_2) \in HK$.

Third, we'll show that HK is closed under inverses. Let $xy \in HK$ with $x \in H$ and $y \in K$. Then $(xy)^{-1} = y^{-1}x^{-1} \in KH = HK$. Q.E.D.

Corollary 4.15. If H and K are subgroups of an Abelian group G , then $H + K$ is also a subgroup of G .

4.2 Symmetric Groups S_n

We've looked at several examples of groups already. It's time to examine some in more detail.

4.2.1 Permutations and the symmetric group

Definition 4.16. A *permutation* of a set X is just a bijection $\rho : X \rightarrow X$ on that set. The permutations on X form a group called the *symmetric group*. We're primarily interested in permutations on a finite set. We'll call the elements of the finite set letters, but we'll denote them with numbers. The symmetric group on n elements $1, 2, \dots, n$ is denoted S_n .

Note that the order of the symmetric group on n letters is $|S_n| = n!$.

Example 4.17. Consider the permutation ρ on set $X = \{1, 2, 3, 4, 5, 6\}$ that exchanges 2 with 4, sends 1 to 3, 3 to 5, and 5 to 1, and leaves 6 fixed. You can describe ρ in a table like this:

n	1	2	3	4	5	6
$\rho(n)$	3	4	5	2	1	6

That table has a lot of redundant information. The first row is just the names of the elements. To describe ρ on an ordered set like X , it's enough to list the elements in the second row: 3, 4, 5, 2, 1, 6. Unfortunately, that makes it harder to figure out where ρ sends an element. The cycle notation, mentioned next, is compact and makes it easier to see how ρ acts. For ρ , this notation will look like $(135)(24)$.

The three elements form a 3-cycle $1 \xrightarrow{\rho} 3 \xrightarrow{\rho} 5 \xrightarrow{\rho} 1$ of ρ denoted (135) . Also note $2 \xrightarrow{\rho} 4 \xrightarrow{\rho} 2$, so (24) is a 2-cycle of ρ . Another name for a 2-cycle is *transposition*. Since $\rho(6) = 6$, therefore (6) by itself is a 1-cycle, also called a *fixed point*, of ρ . The cycle notation for this permutation is $\rho = (135)(24)$. Note that fixed points are not denoted in this notation. Alternatively, this permutation could be denoted $(24)(135)$, or $(531)(42)$, or several other variants.

Since fixed points aren't denoted in cycle notation, we'll need a special notation for the identity permutation since it fixes all points. We'll use 1 to denote the identity since we're using 1 to denote the identity in a group written multiplicatively. In many textbooks the identity is denoted e .

There's a bit of experience needed to quickly multiply two permutations together when they're in cycle notation. Let $\rho = (146)(23)$ and $\sigma = (15)(2643)$. By $\rho\sigma$ mean first perform the permutation ρ then perform σ (in other words, the composition $\sigma \circ \rho$ if we think of these permutations as functions). Then we need simplify the cycle notation

$$\rho\sigma = (146)(23) (15)(2643).$$

Note that first ρ sends 1 to 4, then σ sends 4 to 3, therefore $\rho\sigma$ sends 1 to 3. Next $3 \xrightarrow{\rho} 2 \xrightarrow{\sigma} 6$, so $3 \xrightarrow{\rho\sigma} 6$, likewise $6 \xrightarrow{\rho} 1 \xrightarrow{\sigma} 5$, so $6 \xrightarrow{\rho\sigma} 5$, and $5 \xrightarrow{\rho} 5 \xrightarrow{\sigma} 1$, so $5 \xrightarrow{\rho\sigma} 1$. Thus, we have a cycle of $\rho\sigma$, namely, (1365) . You can check that (2) and (4) are fixed points of $\rho\sigma$. Thus, we found the product. $(146)(23) (15)(2643) = (1365)$.

Incidentally, finding the inverse of a permutation in cycle notation is very easy—just reverse all the cycles. The inverse of $\rho = (146)(23)$ is $\rho^{-1} = (641)(32)$.

Small symmetric groups When $n = 0$ or $n = 1$, there's nothing in the symmetric group except the identity.

The symmetric group on two letters, S_2 , has one nontrivial element, namely, the transposition (12) . This is the smallest nontrivial group, and it's isomorphic to any group of order 2. It is, of course, an Abelian group.

The symmetric group on three letters, S_3 , has order 6. We can name its elements using the cycle notation.

$$1, (12), (13), (23), (123), (132)$$

Besides the identity, there are three transpositions and two 3-cycles. This is not an Abelian group. For instance $(12)(13) = (123)$, but $(13)(12) = (132)$.

The symmetric group on four letters, S_4 , has order 24. Besides the identity, there are $\binom{4}{2} = 6$ transpositions, $\binom{4}{3} \cdot 2 = 8$ 3-cycles, 6 4-cycles, and 3 products of two 2-cycles, like $(12)(34)$.

Exercise 55. Complete the following table listing all 24 of the elements of S_4 .

the identity	1
transpositions	$(12), (13), (14), (23), (24), (34)$
3-cycles	
4-cycles	
products of 2 transpositions	

4.2.2 Even and odd permutations

First we'll note that every cycle, and therefore every permutation, can be expressed as a product of transpositions. We'll soon see after that that a permutation can either be expressed as a product of an even number of transpositions or as a product of an odd number of transpositions, but not both. That will justify the definition of even and odd permutations.

Theorem 4.18. Any cycle can be expressed as a product of transpositions.

Proof. The cycle $(a_1 a_2 a_3 \cdots a_k)$ is the product $(a_1 a_2)(a_1 a_3) \cdots (a_1 a_k)$. Q.E.D.

We'll look at an invariant that will help us distinguish even from odd permutations. It is P_n , the product of all differences of the form $i - j$ where $0 < i < j \leq n$.

$$\begin{aligned}
 P_n &= \prod_{0 < i < j \leq n} (i - j) \\
 &= (1 - 2)(1 - 3) \cdots (1 - n) \\
 &\quad (2 - 3) \cdots (2 - n) \\
 &\quad \dots \\
 &\quad ((n - 1) - n)
 \end{aligned}$$

Lemma 4.19. The effect of applying a transposition to the integers that make up P_n is to change the sign of P_n .

Proof. Let the transposition be (ab) where $0 < a < b \leq n$. The product P_n is made of three factors $P_n = P'P''P'''$ where $P' = (a - b)$, P'' is the product of factors that have either a or b but not both, and P''' is the product of factors that don't have either a or b . Now the transposition (ab) has no effect at all on P''' but negates P' . Its effect on P'' is more complicated. Suppose c is another letter.

Case 1. $c < a < b$. The factors $(c - a)$ and $(c - b)$ of P'' are interchanged by the transposition (ab) .

Case 2. $a < c < b$. The factors $(a - c)$ and $(c - b)$ are interchanged and both negated.

Case 3. $a < b < c$. Like case 1. Thus P'' does not change its value. Since only P' is negated, P_n is negated. Q.E.D.

Theorem 4.20. A permutation is either the product of an even number of transpositions or the product of an odd number of transpositions, but it can't be both.

Proof. Since each transposition negates P_n , the product of an even number of transpositions leaves P_n alone, but the product of an odd number of transpositions negates P_n . It can't be both since P_n is not 0. Q.E.D.

Definition 4.21. A permutation is *even* if it's the product of an even number of transpositions, it's *odd* if it's the product of an odd number of transpositions. The identity 1 is an even permutation.

Note that a cycle is an even permutation if it has an odd length, but it's an odd permutation if it has an even length.

Also, the product of two even permutations is even, the product of two odds is even, and the product of an even and an odd is odd.

Examples 4.22. The symmetric group S_3 has order 6. It's elements are 1, (12), (13), (23), (123), and (132). Three of them, namely 1, (123), and (132) are even while the other three (12), (13), and (23) are odd.

The symmetric group S_4 has 12 even permutations (the identity, eight 3-cycles, and three products of two 2-cycles) and 12 odd permutations (six transpositions and six 4-cycles).

4.2.3 Alternating and dihedral groups

Definition 4.23 (The alternating group A_n). Since the product of even permutations is even, and the inverse of an even permutation is even, therefore the set of even permutations in the symmetric group S_n is a subgroup of S_n . It is called the *alternating group* on n letters, denoted A_n .

For $n \geq 2$, the number of even permutations in S_n is the same as the number of odd permutations, since multiplying by the transposition (12) sets up the bijection. Therefore, the order of A_n is half the order of S_n . So $|A_n| = \frac{1}{2} n!$.

Example 4.24 (Subgroups of S_3). The symmetric group S_3 only has six elements, so it doesn't have many elements. There's the trivial subgroup 1 of order 1. There are three cyclic subgroups of order 2 each isomorphic to C_2 ; besides 1, the other element is one of the transpositions (12), (13) or (23). There's one subgroup of order three, namely, $D_3 = \{1, (123), (132)\}$. (Note that A_3 is the same group as D_3 . The Hasse diagram for the subgroups is fairly simple.

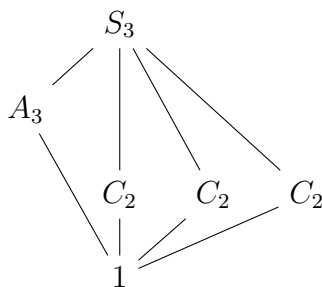


Figure 4.1: Subgroups of S_3

Example 4.25 (The dihedral group D_5). The dihedral groups D_n are the symmetry groups of regular n -gons. We already looked at the case $n = 3$ of an equilateral triangle. Consider a regular polygon with $n = 5$ vertices.

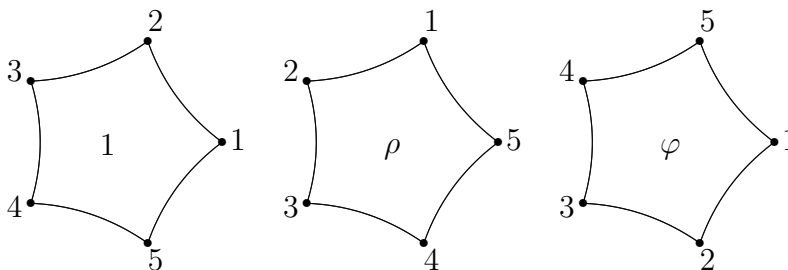


Figure 4.2: Symmetries of a pentagon

We can label the vertices in order from 1 to n . A symmetry of a plane figure is a transformation of the plane that maps the figure to itself. We're only interested in isometries, transformations that preserve distance, right now, but other transformations have their applications, too.

Figure 4.2 shows shows a pentagon. (The pentagon shown here is in the hyperbolic plane, but that doesn't matter.) One of its symmetries ρ is the one that rotates the pentagon 72° counterclockwise. It maps the vertex labelled 1 to 2, maps 2 to 3, and so forth. Knowing where the vertices are mapped is enough to determine the transformation, so we can identify ρ with the permutation it describes on the set of vertices. This ρ is the permutation (12345) .

Another of the symmetries of the pentagon is a reflection like φ shown above, a reflection across a horizontal axis. In cycle notation $\varphi = (25)(34)$.

In fact, there are 10 symmetries of the regular pentagon, so $|D_5| = 10$. In general $|D_n| = 2n$. In D_5 , besides the identity, there are four rotations and five reflections.

$$\begin{array}{lllll} \text{identity} = 1 & \rho = (12345) & \rho^2 = (13524) & \rho^3 = (14253) & \rho^4 = (15432) \\ \varphi = (25)(34) & \varphi\rho = (12)(35) & \varphi\rho^2 = (13)(45) & \varphi\rho^3 = (14)(23) & \varphi\rho^4 = (15)(24) \end{array}$$

There are no more symmetries although we can write more expressions in terms of φ and ρ , for instance $\rho\varphi$. But $\rho\varphi = (15)(24)$ which is $\varphi\rho^4$.

Thus, we can see now how to represent the dihedral group, D_5 , as a subgroup of the symmetric group S_5 . In fact, it's represented as a subgroup of the alternating group, A_5 as well, since all the permutations are even permutations.

Example 4.26 (Symmetries of a cube and tetrahedron). Consider a cube with vertices $12'34'1'2'3'4'$ and the inscribed regular tetrahedron 1234 shown in figure 4.3. The four diagonals of the cube, $11'$, $22'$, $33'$, and $44'$, are drawn in green.

There are many symmetries of a tetrahedron. They permute the vertices 1234 . There are rotations of 120° and 240° about any of the four diagonals. Those rotations about the line $11'$ are the permutations (234) and (243) . The rotations about the other three diagonals are (123) , (132) , (124) , (142) , (134) , and (143) . Besides these rotations, there are three 180° rotations about the three lines joining the midpoints of the opposite edges of the tetrahedron. Along with the identity, that makes 12 permutations, all of which preserve orientation, that

is, they're rigid motions. The group of orientation preserving symmetries of the tetrahedron form the group S_4 .

Besides these, there are symmetries of the tetrahedron which are reflections across planes. They are orientation reversing symmetries. For example, the reflection across the plane passing through vertices 1 and 2 and the midpoint of edge 34 leaves vertices 1 and 2 fixed but it exchanges vertices 3 and 4; it's the transposition (34) . The group of all the symmetries of the tetrahedron, including both the orientation preserving symmetries and the orientation reversing ones, form the group S_4 .

Each of the symmetries of the tetrahedron gives a symmetry of the enclosing cube. The symmetries of a cube permute its eight vertices $12341'2'3'4'$. For example, the symmetry (123) of the tetrahedron gives the symmetry $(123)(1'2'3')$ of the cube.

But there are other symmetries of a cube since the tetrahedron 1234 doesn't have to be preserved under a symmetry of the cube; it could be sent to the opposite tetrahedron $1'2'3'4'$. Other orientation preserving symmetries that send the tetrahedron to the opposite tetrahedron include the six 90° and 270° rotations about the centers of the faces and the six 180° rotations about the line joining midpoints of opposite sides. That makes 24 orientation preserving symmetries for the cube. Each one permutes the four diagonals, and no two of them permute the four diagonals in the same way, so this symmetry group is S_4 .

Note that the symmetry $(11')(22')(33')(44')$ that exchanges a vertex with its opposite vertex reverses orientation.

The entire group of symmetries of the cube includes the 24 orientation preserving symmetries and each of those times $(11')(22')(33')(44')$. That makes 48 symmetries of the cube.

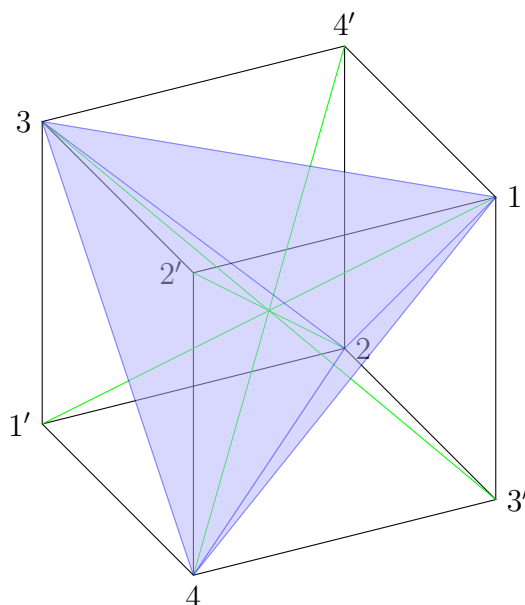


Figure 4.3: Symmetries of a cube and tetrahedron

Exercise 56. Verify the statements made in the example.

- The orientation preserving symmetries of a tetrahedron form the group A_4 .
- The group of all the symmetries of the tetrahedron form the group S_4 .

- (c). The orientation preserving symmetries of a tetrahedron form the group S_4 .
 (c). Explain why the permutation $(11')(22')(33')(44')$ of the cube reverses orientation.

Presentations by generators and relations. Although it's nice to have a group represented in a symmetric group, sometimes it's more convenient to describe it more algebraically in terms of generators and relations. For D_5 we can see that ρ and φ are sufficient to generate the whole group in the sense that every element in the group can be written as some expression involving ρ and φ . But there are certain relations, actually equations, that ρ and φ satisfy in this group, namely $\rho^5 = 1$, $\varphi^2 = 1$, and $\rho\varphi = \varphi\rho^{-1}$. Thus, we can present the group as

$$D_5 = \langle \rho, \varphi : \rho^5 = 1, \varphi^2 = 1, \rho\varphi = \varphi\rho^{-1} \rangle.$$

The difficulty with a presentation of this type is knowing when you have enough generators and relations. If you don't have enough generators, you won't generate the whole group. If you don't have enough relations, you'll generate a larger group, but not the one you want. A proof needs to be supplied to be assured that this is the right presentation. Frequently, a diagram of some sort fills the bill.

4.3 Cayley's theorem and Cayley graphs

One of the reasons symmetric groups are so important is that every group is isomorphic to a subgroup of a symmetric group, a result of Cayley. This gives us another way to look at groups, especially small finite ones.

We'll prove Cayley's theorem, then look at a few Cayley graphs which depend on Cayley's theorem.

4.3.1 Cayley's theorem

Recall that a permutation of a set X is just a bijection $\rho : X \rightarrow X$ on that set and permutations on X form a group called the *symmetric group* $S(X)$. When the set is finite, we can write it as $\{1, 2, \dots, n\}$, and S_n denotes the its symmetric group.

Cayley's theorem can be stated for infinite groups as well as finite groups.

Theorem 4.27 (Cayley). Let G be a group, and let $S(G)$ be the symmetric group on G , that is, the group of permutations on the underlying set of G . The function $\varphi : G \rightarrow S(G)$ defined by $\varphi(a)(x) = ax$ is a group monomorphism. Therefore, G is isomorphic to a subgroup of $S(G)$.

Proof. $\varphi(a)$ is the permutation on G that maps x to ax . It's a bijection since its inverse sends x to $a^{-1}x$. To show that it's a group homomorphism, it is only necessary to show that $\varphi(ab) = \varphi(a)\varphi(b)$ for a and b in G . But $\varphi(ab)(x) = abx$, and $(\varphi(a)\varphi(b))(x) = \varphi(a)(\varphi(b)(x)) = \varphi(a)(bx) = abx$. Finally, $\varphi : G \rightarrow S(G)$ is a monomorphism since if $\varphi(a) = \varphi(b)$, then evaluating the two permutations at 1 gives $a1 = b1$, so $a = b$. Q.E.D.

Although this representation theorem does show that every group is a subgroup of a symmetric group (up to isomorphism), it's practically not all that useful since if the group G has order n , it's being represented in a group of order $n!$, which is much too large to deal with if n is at all large. Still, it's a useful representation for theoretical purposes.

Cayley graphs. With a Cayley graph we can represent a group G by a graph with vertices and labeled, directed edges. Each element of G is a vertex of the graph, and for each element a , we also have a directed edge labeled a from a vertex x to the vertex ax . In other words, the Cayley graph is a representation of G by the Cayley theorem to $S(G)$.

For a small example, let G be the cyclic group $G = \{1, a, b\}$ where $a^2 = b$ and $a^3 = 1$. The Cayley graph for G has three vertices, labeled $1, a,$ and b . Each node has a loop on it labeled 1 since $1x = x$. There are three edges labelled a , $1 \xrightarrow{a} a \xrightarrow{a} b \xrightarrow{a} 1$, and three edges labelled b , $1 \xrightarrow{b} b \xrightarrow{b} a \xrightarrow{b} 1$. This is probably most conveniently drawn in a triangular figure.

There's a lot of redundancy in the graph in the sense that you don't need all the information to reconstruct the group. The loops labelled 1 might just as well be dropped since for any group $1x = x$. If we know the edges labelled a , then we can determine the edges labelled b since you just travel two a -edges to get a b -edge. That leaves just the triangle $1 \xrightarrow{a} a \xrightarrow{a} b \xrightarrow{a} 1$. More generally, if you know the edges for generators of a group, then all the other edges are determined.

Example 4.28 (D_5). Recall that the dihedral group D_5 has 10 elements and the presentation

$$D_5 = \langle \rho, \varphi : \rho^5 = \varphi^2 = (\varphi\rho)^2 = 1 \rangle.$$

The first relation, $\rho^5 = 1$ gives us a five cycle

$$1 \xrightarrow{\rho} \rho \xrightarrow{\rho} \rho^2 \xrightarrow{\rho} \rho^3 \xrightarrow{\rho} \rho^4 \xrightarrow{\rho} 1$$

which we can draw as a pentagon, the center pentagon in the graph below. The second relation, $\varphi^2 = 1$, means we have the 2-cycle $1 \xrightarrow{\varphi} \varphi \xrightarrow{\varphi} 1$, and, more generally, for any element a , we have a 2-cycle $a \xrightarrow{\varphi} a\varphi \xrightarrow{\varphi} a$. We'll draw 2-cycles as undirected edges $a \overset{\varphi}{\dashrightarrow} a$. We get five of these edges, one at each vertex of the center pentagon. The third relation, $(\varphi\rho)^2 = 1$, describes a square

$$a \overset{\varphi}{\dashrightarrow} a\varphi \xrightarrow{\rho} a\varphi\rho \overset{\varphi}{\dashrightarrow} a\varphi\rho\varphi \xrightarrow{\rho} a.$$

Starting at each of the new outer vertices of the graph, follow three edges to reach another outer vertex, and draw a ρ -edge back to where you started. When you finish, you have the Cayley graph for D_5 in figure 4.4

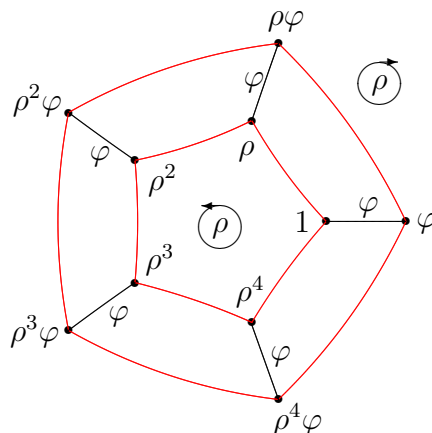


Figure 4.4: Cayley graph for D_5

Notice that the graph is completely symmetric. You could label any vertex 1 and fill in the names of the rest of the vertices by following the labelled arcs. For that reason, the vertices of a Cayley graph needn't be labelled.

There is another presentation for D_5 that gives a different looking Cayley graph. Let $\psi = \rho\phi$. Then

$$D_5 = \langle \varphi, \psi : \varphi = \psi^2 = (\varphi\psi)^5 \rangle.$$

The Cayley graph has the same ten vertices, but the edges are all undirected and they form a cycle of length 10 with labels alternating between φ and ψ .

Example 4.29 (A_4). Recall that the alternating group on $\{1, 2, 3, 4\}$ has 12 elements. It's not cyclic, so at least two generators are required to generate it. In fact, two will do. Consider the three elements

$$\begin{aligned} a &= (123) \\ b &= (124) \\ c &= ab = (14)(23) \end{aligned}$$

The two elements a and b are sufficient to generate A_4 as are the two elements a and c and many other pairs of elements (but not all pairs will do). In fact, A_4 can be represented in either of the following two ways:

$$\begin{aligned} \langle a, b : a^3 = b^3 = (ab)^2 = 1 \rangle \\ \langle a, c : a^3 = c^2 = (ac)^2 = 1 \rangle \end{aligned}$$

So, if we have the Cayley graph with only a - and b -edges, then we have enough information to determine A_4 , or if we have the graph with only a - and c -edges, then that's enough. Although these two graphs both have 12 vertices (since $|A_4| = 12$), they don't look very much alike. Let's look at the Cayley graph with all three kinds of edges, a -edges and b -edges and c -edges.

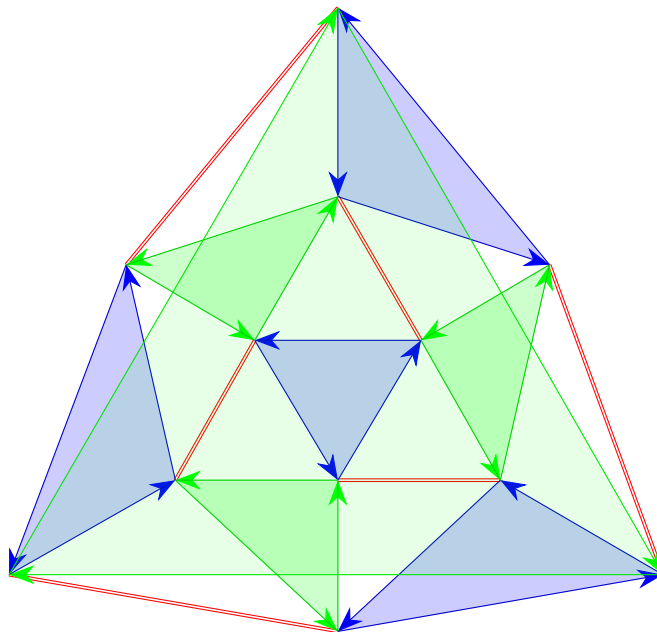
It's displayed in figure 4.5 as a planar graph, but more of the symmetry would be apparent if it were displayed in three dimensions where the vertices and edges were those of an icosahedron. Some of the triangles in the figure are blue. Their three sides of the triangle are a -edges. Likewise, some triangles are green with b -edges. Note that all the a - and b -triangles are oriented counterclockwise except the outer b -triangle. The remaining edges are the red c -edges, and to save space, since c is an involution, rather than putting in two edges, one pointing one way and the other pointing the other way, just a single thick undirected edge is included. Each vertex in the graph has an a -edge coming in and one coming out, a b -edge coming in and one coming out, and an undirected c -edge meaning that it goes both in and out.

Since it only takes two of these three elements to generate A_4 , this graph has superfluous information. All the edges labelled by one of the letters can be removed making the graph simpler.

Exercise 57. Find a Cayley graph for the symmetric group S_4 . There are various pairs or triples of generators you can use. One is the pair $a = (1234), b = (12)$.

4.3.2 Some small finite groups

We've seen a few families of finite groups including C_n the cyclic group of order n , D_n the dihedral group of order $2n$, S_n the symmetric group of order $n!$, and A_n the alternating group of order $n!/2$.

Figure 4.5: Cayley graph for A_4

The classification of finite groups (up to isomorphism, of course) is extremely difficult. Daniel Gorenstein (1923–1992) was a leader of mathematicians who eventually classified finite simple groups. He was faculty member at Clark University for 13 years.

We'll look at a few more small finite groups. Later, we'll look at the classification of finite Abelian groups, and find that they're all products of cyclic groups.

Table 4.1 lists the small groups up to isomorphism of order up through 24.

order	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
number of groups	1	1	1	2	1	2	1	5	2	2	1	5	1	2	1

We won't prove that these are all of them, but we will look at them all. There are combinatorial theorems, the most important being the Sylow theorems, that help in classifying finite groups.

We know nearly all of these 27 groups. The cyclic groups C_n account for 15 of them. There are 12 others. Some of them are products of smaller ones, for instance, the other group of order 4 is $C_2 \oplus C_2$, sometimes called the *Klein 4-group*.

The second group of order 6 is D_3 , which is the same as S_3 .

Two of the groups of order 8 are products, namely, $C_4 \oplus C_2$ and $C_2 \oplus C_2 \oplus C_2$. Another is D_4 and the remaining one is called the quaternion group.

Example 4.30 (The quaternion group). This group consists of eight of the units of the division ring \mathbf{H} , the quaternions. Let $Q = \{\pm 1, \pm i, \pm j, \pm k\}$. Recall that the multiplication of quaternions has $i^2 = j^2 = k^2 = -1$, $ij = k$, $jk = i$, and $ki = j$, so this set of units is closed under multiplication and forms a group, called the *quaternion group*.

Exercise 58. Construct a Cayley graph for the quaternion group.

The second group of order 9 is the Abelian group $C_3 \oplus C_3$, and the second group of order 10 is D_5 . We already know the other groups of order 12: D_6 , $C_2 \oplus C_6$, $D_3 \times C_2$, and A_4 , and

Order	1	2	3	4	5	6	7	8
Abelian groups	C_1	C_2	C_3	C_4 $C_2 \oplus C_2$	C_5	C_6	C_7	C_8 $C_4 \oplus C_2$ $C_2 \oplus C_2 \oplus C_2$
Non-Abelian						D_3		D_4 Q
Order	9	10	11	12	13	14	15	16
Abelian groups	C_9 $C_3 \oplus C_3$	C_{10}	C_{11}	C_{12} $C_2 \oplus C_6$	C_{13}	C_{14}	C_{15}	C_{16} $C_2 \oplus C_8$ $C_2 \oplus C_2 \oplus C_4$ $C_2 \oplus C_2 \oplus C_2 \oplus C_2$ $C_4 \oplus C_4$
Non-Abelian		D_5		D_6 A_4 $C_3 \times C_4$		D_7		D_8 8 others
Order	17	18	19	20	21	22	23	24
Abelian groups	C_{17}	C_{18} $C_3 \oplus C_6$	C_{19}	C_{20} $C_2 \oplus C_{10}$	C_{21}	C_{22}	C_{23}	C_{24} $C_2 \oplus C_{12}$ $C_2 \oplus C_2 \oplus C_6$
Non-Abelian		D_9 $S_3 \times C_3$ $C_3^2 \rtimes C_2$		D_{10} Dic_5 $C_5 \rtimes C_4$	$C_7 \rtimes C_3$	D_{11}		D_{12} 11 others

Table 4.1: List of small groups

the other group of order 14 is D_7 .

4.4 The category of groups \mathcal{G}

The category \mathcal{G} of groups was mentioned briefly in section 3.5 when the category of rings was introduced. The objects in this category are groups, and the morphisms are group homomorphisms.

Products in a category are defined by a universal property rather than by ordered pairs, but in \mathcal{G} , the product of two groups is what was called the the product of groups in section 4.1.4:

$$G \times H = \{(x, y) \mid x \in G, y \in H\}$$

Other categorical concepts include the initial and final object. In \mathcal{G} these are the same group, namely the trivial group that has only one element 1. There is a unique group homomorphism from each group to 1, and there's a unique group homomorphism from 1 to each group.

The universal property of an infinite cyclic group in the category of groups. The addition operation on \mathbf{Z} makes it an infinite cyclic group since each element in it is a multiple of 1. You can also write it multiplicatively as $C_\infty = \langle a \rangle = \{\dots, a^{-2}, a^{-1}, 1, a, a^2, \dots\}$.

This infinite cyclic group has the following universal property. Given any group G and any element $c \in G$, there is a unique group homomorphism $\langle a \rangle \rightarrow G$ that maps a to c . In general, it maps a^n to c^n . The image of this homomorphism is the subgroup of G generated by c .

Free groups. The infinite cyclic group $\langle a \rangle$ is a special case of a free group. It's a free group on one element. There are free groups on more than one element with analogous universal properties. We'll look at the free group on two elements.

Let a and b be two symbols. Form the group $\langle a, b \rangle$ as follows. An element in it is named by a string of a 's and b 's raised to various integral powers, such as $ba^{-3}b^2b^4a^2$. Different names are to be considered to name the same element if adjacent symbols are the same, in which case, they can be combined by the usual power rule. For example, $ba^{-3}b^2b^4a^2 = ba^{-3}b^6a^2$. Also, any symbol to the power 0 is to be treated as the identity element 1, and 1 times any string simplifies to that string.

A formal proof that $\langle a, b \rangle$ is, in fact, a group requires induction. We'll omit that proof.

This group $\langle a, b \rangle$ has the following universal property. Given any group G and any elements $c, d \in G$, there is a unique group homomorphism $\langle a, b \rangle \rightarrow G$ that maps a to c and b to d .

4.5 Conjugacy classes and quandles

We'll consider another way to examine the structure of groups. That depends on analyzing the operation of conjugation in a group.

Definition 4.31 (Conjugate element in a group). If x and y are elements of a group G , then xyx^{-1} is called *conjugates* of x . In that case, elements y and xyx^{-1} are said to be conjugates.

Exercise 59. Show that being an conjugate in a group is an equivalence relation by proving that (a) any element is conjugate to itself, (b) if one element is conjugate to second, then the second is conjugate to the first, and (c) if one element is conjugate to a second and the second conjugate to a third, then the first is conjugate to the third.

4.5.1 Conjugacy classes

Since being conjugates in a group is an equivalence relation, the corresponding equivalence classes can say a lot about the group.

Definition 4.32. Each of the equivalence classes of a group under conjugacy is called a *conjugacy class*, and the set of all conjugates of a particular element x is called the *conjugacy class* of x .

Exercise 60. If x is an element of order n in a group, show every conjugate of x also has order n .

Example 4.33 (Conjugacy classes in symmetric groups). Conjugation and conjugacy classes in symmetric groups are particularly easy to identify using cycle notation. Let $x = (13)(245)$ and $y = (142)$ be two elements in S_n . Then $y^{-1}xy = (124)(13)(245)(142) = (43)(125)$. Note how y conjugates the cycle (13) to the cycle (43) , and it conjugates the cycle (245) to (125) . The cycle structures for x and $y^{-1}xy$ are the same, but the elements in the cycles are permuted by y . This is generally the case for symmetric groups. It follows that a conjugacy class in S_n consists of all the elements in S_n with a given structure. Thus, for example, the conjugacy class of $(13)(235)$ consists of all elements of the form $(ab)(cde)$ where a, b, c, d , and e are 5 distinct integers between 1 and n . For S_5 the size of that conjugacy class is $\binom{5}{2} \cdot 2 = 20$.

Exercise 61. Determine all the conjugacy classes of S_5 and their sizes. (The sum of their sizes will equal 120, of course.)

Theorem 4.34. If H is a subgroup of G , and $x \in G$, then xHx^{-1} is also a subgroup of G , called a subgroup *conjugate* to H .

Proof. First, $1 \in xHx^{-1}$ since $x1x^{-1} = 1$. Next, if $xyx^{-1}, xzx^{-1} \in xHx^{-1}$ with $y, z \in H$, then their product $xyx^{-1}xzx^{-1} = x(yz)x^{-1} \in xHx^{-1}$. Finally, given $xyx^{-1} \in xHx^{-1}$ with $y \in H$, then the inverse $(xyx^{-1})^{-1} = xy^{-1}x^{-1} \in xHx^{-1}$. Therefore, xHx^{-1} is a subgroup of G . Q.E.D.

Similarly to the argument in the exercise above, being conjugate subgroups of a given group is an equivalence relation.

Theorem 4.35. If no other subgroup of G has the same order as H , then H is normal.

Proof. Since any conjugate subgroup xHx^{-1} is in one-to-one correspondence with H , it has the same number of elements, so must equal H . Q.E.D.

Exercise 62. If H is a subgroup of G and N is a normal subgroup of G , prove that $H \cap N$ is a normal subgroup of H .

Exercise 63. If H is a subgroup of G and N is a normal subgroup of G , prove that HN is a subgroup of G . (Hint: show $HN = NH$.)

Exercise 64. Prove that the intersection of two normal subgroups is also a normal subgroup.

Exercise 65. Prove that if H and N are normal subgroups of G , then their product is also a normal subgroup of G , in fact, it's the subgroup generated by $H \cup N$.

4.5.2 Quandles and the operation of conjugation

The operations of conjugation have certain properties. If we think of $y^{-1}xy$ as a binary operation $x \triangleright y$, and xyx^{-1} as another operation $x \triangleright^{-1} y$, then these two operations satisfy the properties stated in the next definition.

Definition 4.36. A *quandle* is a set equipped with two operations, \triangleright and \triangleright^{-1} satisfying the following three conditions for all elements x, y , and z .

- Q1.** $x \triangleright x = x$.
- Q2.** $(x \triangleright y) \triangleright^{-1} y = x = (x \triangleright^{-1} y) \triangleright y$.
- Q3.** $(x \triangleright y) \triangleright z = (x \triangleright z) \triangleright (y \triangleright z)$.

The symbol \triangleright is pronounced “through”, and \triangleright^{-1} “backthrough”.

Exercise 66. Prove that if Q is a conjugacy class in a group G then Q is a quandle where the operation $x \triangleright y$ is ${}^{-1}xy$, and $x \triangleright^{-1} y$ is xyx^{-1} .

Involutory quandles. A quandle satisfying the identity $x \triangleright y \triangleright y = x$, equivalently $x \triangleright y = x \triangleright^{-1} y$, is called an *involutory quandle* or a *2-quandle*. The two operations of a quandle are the same in an involutory quandle.

There is an analogous definition for an n -quandle. First define $x \triangleright^n y$ as $x \triangleright y \triangleright \cdots \triangleright y$ where $\triangleright y$ occurs n times. An n -quandle is a quandle that satisfies the identity $x \triangleright^n y = x$.

A conjugacy class of involutions in a group is an involutory quandle, while the conjugacy class of an element of order n is an n -quandle.

Conjugacy classes of involutions are useful in the study of groups.

Besides conjugacy classes of groups, involutory quandles appear as cores of a group. The *core* of a group G has the same elements as G but with the operation $x \triangleright y = yx^{-1}y$.

Exercise 67. Prove that the core of a group is an involutory quandle.

Involutory quandles with geodesics. Involutory quandles have a nice geometric interpretation where the elements are points and the lines are determined by the operation.

Example 4.37 (The plane as a quandle). Consider the Euclidean plane \mathbf{R}^2 with the operation which sends a point p through a point q to yield the point $p \triangleright q$ on the line that passes through p and q and on the opposite side of q that p lies on but equally far away from q . If p happens to equal q , then define $p \triangleright q$ to be q . Algebraically, $p \triangleright q = 2q - p$.

This operation makes \mathbf{R}^2 an involutory quandle. The self distributivity axiom **Q3**, which says $(p \triangleright q) \triangleright r = (p \triangleright r) \triangleright (q \triangleright r)$, is illustrated in figure 4.6.

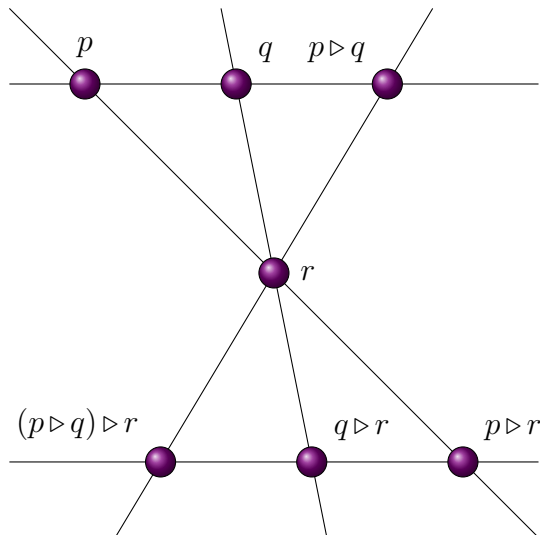


Figure 4.6: Distributivity in an involutory quandle

Symmetric spaces. A symmetric space is a particular kind of manifold. At each point in the space, there is an isometry (that is, a translation that preserves distance) for which that point is an isolated singularity. Ottmar Loos discovered in 1967 that the intrinsic algebraic structure of a symmetric is an involutory quandle. Thus, a symmetric space as a differentiable involutory quandle in which every point is an isolated fixed point of the symmetry through it.

Besides the plane \mathbf{R}^2 , every vector spaces V over any field is a symmetric spaces. The operation that makes it an involutory quandle is given by $\mathbf{v} \triangleright \mathbf{w} = 2\mathbf{w} - \mathbf{v}$. That's the same operation as described above for \mathbf{R}^2 .

There are lots of other symmetric spaces. The ordinary sphere S^2 as well as higher dimensional spheres S^n are all symmetric spaces. So are other geometric spaces including projective spaces, hyperbolic spaces, and inversive spaces. They can all be used to represent quandles geometrically as subspaces.

Geodesics. A *geodesic* in a manifold is a curve which for points close together is the curve of shortest length that joins them. In Euclidean space, a geodesic is a straight line. On the sphere S^2 , a geodesic is a great circle, that is, the intersection of a plane passing through the center of the sphere with the sphere. Geodesics on manifolds have metrics, that is, there's a distance between any two points on the geodesic.

Given two points p and q , the entire involutive quandle generated by them lies on one geodesic. That means that any other expression that can be made from p and q lie on a geodesic. In particular, the points $p \triangleright (q \triangleright p)$, $q \triangleright p$, p , q , $p \triangleright q$, and $q \triangleright (p \triangleright q)$ lie on a geodesic, and they're equally spaced on it.

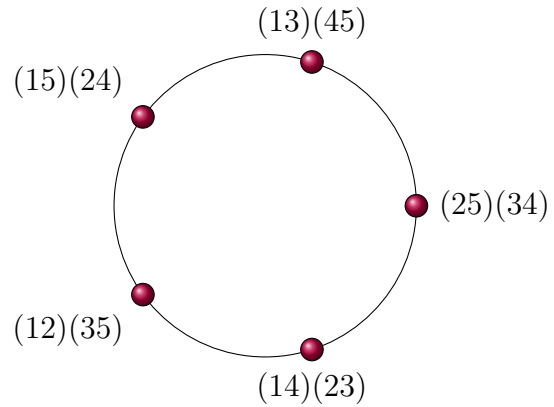


Since geodesics have all the information needed to describe the quandle structure, we can draw involutory quandles, at least the small ones, to see what they look like. Some will be

generated by two elements, so they'll either look like the line above, or be quotients of it.

Example 4.38 (Reflections of a pentagon). The group D_5 of symmetries of the pentagon was illustrated in figure 4.2.

The reflections of a pentagon are involutions, and they form a conjugacy class in D_5 . There are five reflections: $p = (25)(34)$, $q = (13)(25)$, $p \triangleright q = (15)(24)$, $q \triangleright p = (14)(23)$, and $p \triangleright (q \triangleright p) = q \triangleright (p \triangleright q) = (12)(35)$. Since $p \triangleright (q \triangleright p) = q \triangleright (p \triangleright q)$, the five of them lie on a circle as illustrated to the right.



Exercise 68. Let the seven vertices of a regular heptagon be denoted 1, 2, 3, 4, 5, 6, and 7. Describe the how the symmetry $(27)(36)(45)$ acts on the heptagon in words. Determine the conjugacy class of $(27)(36)(45)$ in D_7 and illustrate it as points on a circle.

Some conjugacy classes of involutions are cyclic like the ones in D_5 above, but most aren't. Here are two examples of 6-element conjugacy classes in small groups.

Example 4.39 (A conjugacy class in the quaternion group). The quaternion group was introduced in section 4.30. It has eight elements, namely 1, -1 , i , $-i$, j , $-j$, k , and $-k$. Six of them, all those except ± 1 , are involutions and they form a conjugacy class. It's illustrated in figure 4.7. Note that $i \triangleright j = j i j = -i$, so i , j , $-i$, and $-j$ are equally spaced around a circle. Likewise for i , k , $-i$, and $-k$ and for k , j , $-k$, and $-j$. Although the spacing doesn't appear equal on the Euclidean plane as shown, it is when represented on a sphere.

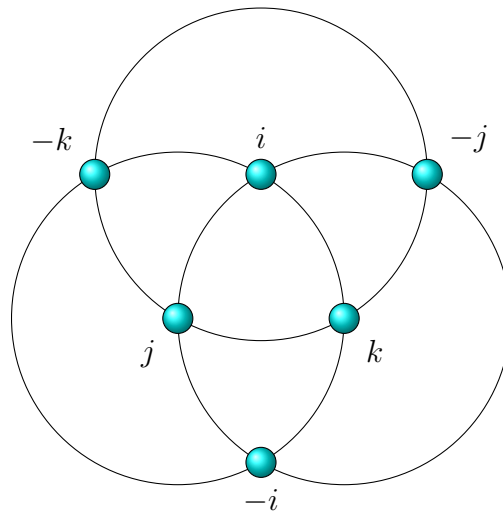


Figure 4.7: A conjugacy class in the quaternion group

Example 4.40 (The conjugacy class of transpositions in S_4). The transpositions in a symmetric group form a conjugacy class. The the symmetric group S_4 there are six transpositions,

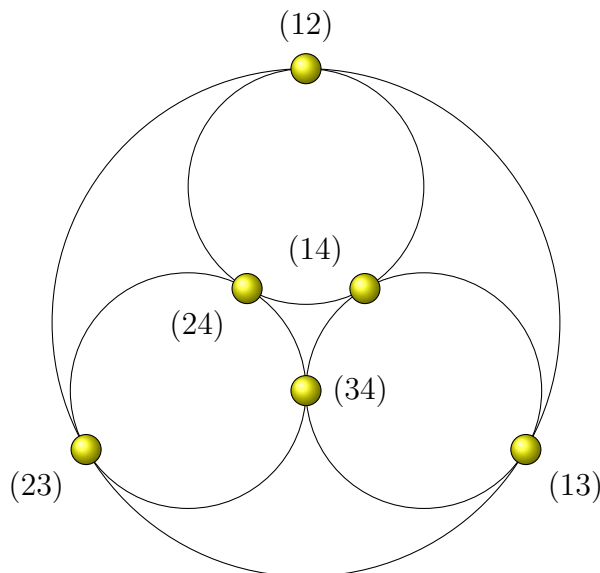


Figure 4.8: The conjugacy class of transpositions in S_4

namely, (12), (13), (14), (23), (24), and (34). The involutory quandle they form is shown in figure 4.8. There are four circles, each with three transpositions. One circle, for example, includes (12), (13), and (23) since $(12) \triangleright (13) = (23)$, $(13) \triangleright (23) = (12)$, and $(23) \triangleright (12) = (13)$. Also note that $(12) \triangleright (34) = (12)$, but no geodesic is shown having those two transpositions. It reduces the clutter in the diagram to suppress geodesics with only two elements.

4.6 Kernels, normal subgroups, and quotient groups

The kernel $\text{Ker } f$ of a group homomorphism $f : G \rightarrow H$ plays the same role as the kernel of a ring homomorphism. It's defined as the the inverse image of the identity. It is a subgroup of the domain G , but a particular kind of subgroup called a normal subgroup. We'll see that every normal subgroup N of G is the kernel of some group homomorphism, in fact, of a projection $G \rightarrow G/N$ where G/N is a quotient group of G .

4.6.1 Kernels of group homomorphisms and normal subgroups

We'll use multiplicative notation.

Definition 4.41. Let $f : G \rightarrow H$ be a group homomorphism. Those elements of G that are sent to the identity 1 in H form the *kernel* of f .

$$\text{Ker } f = f^{-1}(1) = \{x \in G \mid f(x) = 1\}.$$

Example 4.42. Let G be the symmetric group S_n and $f : G \rightarrow \{1, -1\}$ map even permutations to 1 and odd permutations to -1 . Then f is a group homomorphism, and $\text{Ker } f = A_n$, the alternating subgroup of S_n .

Theorem 4.43. The kernel of a group homomorphism $f : G \rightarrow H$ is a subgroup $N = \text{Ker } f$ of G such that for each $x \in G$, $xNx^{-1} \subseteq N$.

Proof. To show that N is a subgroup of G , note that (1) it's closed under multiplication, (2) it includes 1, and (3) it's closed under inverses. For (1), if $x, y \in N$, then $f(x) = f(y) = 1$, so $f(xy) = f(x)f(y) = 1$, therefore $xy \in N$. (2) is obvious. For (3), if $x \in N$, then $f(x) = 1$, so $f(x^{-1}) = f(x)^{-1} = 1^{-1} = 1$, therefore $x^{-1} \in N$. Thus N is a subgroup of G .

Now to show that for $x \in G$, $xNx^{-1} \subseteq N$. Consider xyx^{-1} where $y \in N$. Then $f(y) = 1$, so $f(xyx^{-1}) = f(x)f(y)f(x)^{-1} = f(x)1f(x)^{-1} = f(x)f(x)^{-1} = 1$. Therefore, $xyx^{-1} \in N$. Thus, $xNx^{-1} \subseteq N$. Q.E.D.

Besides telling us what elements are sent to 1 by f , the kernel of f also tells us when two elements are sent to the same element. Since $f(x) = f(y)$ if and only if $f(xy^{-1}) = 1$, therefore, f will send x and y to the same element of S if and only if $xy^{-1} \in \text{Ker } f$.

The properties of kernels of group homomorphisms that we just found determine the following definition.

Definition 4.44. A subgroup N of a group G is said to be a *normal subgroup* if for each $x \in G$, $xNx^{-1} \subseteq N$.

Note that since a normal subgroup a group G is closed under conjugation, therefore a normal subgroup of G is the union of some of the conjugacy classes in G .

Exercise 69. Show that a subgroup N is normal in G if and only if for each $x \in G$, $xNx^{-1} = N$.

Exercise 70. Show that a subgroup N is normal in G if and only if for each $x \in G$, $xN = Nx$.

Both the trivial subgroup of G and G itself are always normal subgroups.

If G is an Abelian group, then every subgroup of G is a normal subgroup.

Theorem 4.45. Any subgroup of index 2 is a normal subgroup.

Proof. Let N be a subgroup of a group G of index 2. We'll show that $xN = Nx$ for each $x \in G$. In case $x \in N$, then $xN = N = Nx$. Now consider the case $x \notin N$. Then there are two left cosets of N , namely N itself and xN , and there are two right cosets, N and Nx . That gives us two partitions of G , but since N is a part of each partition, the other parts, namely xN and Nx must be equal. Q.E.D.

4.6.2 Quotient groups, and projections $\gamma : G \rightarrow G/N$

As mentioned above the kernel of a group homomorphism f tells us when two elements are sent to the same element: $f(x) = f(y)$ if and only if $xy^{-1} \in \text{Ker } f$. We can use $\text{Ker } f$ to construct a "quotient group" $G/\text{Ker } f$ by identifying two elements x and y in G if xy^{-1} lies in $\text{Ker } f$. In fact, we can do this not just for kernels of homomorphisms, but for any normal subgroup N . That is, we can use a normal subgroup N of G to determine when two elements x and y are to be identified, $x \equiv y$, and we'll end up with a group G/N .

Definition 4.46. A *congruence* \equiv on a group G is an equivalence relation such that for all $x, x', y, y' \in G$,

$$x \equiv x' \text{ and } y \equiv y' \text{ imply } xy \equiv x'y'.$$

The equivalence classes for a congruence are called *congruence classes*.

Theorem 4.47. If \equiv is a congruence on a group G , then the quotient set G/\equiv , that is, the set of congruence classes, is a group where the binary operation is defined by $[x][y] = [xy]$.

Proof. First we need to show that the proposed definitions are actually well defined. That is, if a different representative x' is chosen from the congruence class $[x]$ and y' from $[y]$, then the same class $[x'y']$ results. That is

$$[x] = [x'] \text{ and } [y] = [y'] \text{ imply } [xy = xy'].$$

But that is the requirement in the definition of congruence.

Also, each of the axioms for a group need to be verified, but they're all automatic as they're inherited from the group G . Q.E.D.

Just as an ideal in a ring determines a congruence on the ring, a normal subgroup of a group determines a congruence on a group, and the proof is similar.

Theorem 4.48 (Congruence modulo a normal subgroup). Let N be a normal subgroup of a group G . A congruence, called *congruence modulo N* , is defined by

$$x \equiv y \pmod{N} \text{ if and only if } xy^{-1} \in N.$$

The quotient group, G/\equiv , is denoted G/N . The congruence classes are cosets of N , that is $[x] = xN$. The function $\gamma : G \rightarrow G/N$ defined by $\gamma(x) = [x] = xN$ is a group homomorphism, in fact, an epimorphism. It's called a *projection* or a *canonical homomorphism* to the quotient group. Its kernel is N .

Exercise 71. If \equiv is a congruence on a group G , show that the congruence class of the identity, $[1] = N$, is a normal subgroup of G , and the congruence determined by N is the original congruence.

4.6.3 Isomorphism theorems

The image of a group homomorphism is isomorphic to the group modulo its kernel. Let $f : G \rightarrow H$ be a ring homomorphism. The image of f , denoted $f(G)$, is the set

$$f(G) = \{f(x) \in H \mid x \in G\}.$$

Exercise 72. Verify that the image $f(G)$ is a subgroup of H .

Exercise 73. Prove the following theorem. You'll need to show that the proposed function is well-defined, that it is a group homomorphism, and then that it's an isomorphism.

Theorem 4.49 (First isomorphism theorem, Jordan, 1870). If $f : G \rightarrow H$ is a group homomorphism then the quotient group $G/\text{Ker } f$ is isomorphic to the image ring $f(G)$, the isomorphism being given by

$$\begin{aligned} G/\text{Ker } f &\rightarrow f(G) \\ x \text{Ker } f &\mapsto f(x) \end{aligned}$$

This gives us two ways to look at the image, either as a quotient group of the domain G or as a subgroup of the codomain H .

Furthermore, we can now treat a group homomorphism $f : G \rightarrow H$ as a composition of three group homomorphisms.

$$G \xrightarrow{\gamma} G/\text{Ker } f \cong f(G) \xrightarrow{\iota} H$$

The first is the projection from G onto its quotient ring $G/\text{Ker } f$, the second is the isomorphism $G/\text{Ker } f \cong f(G)$, and the third is the inclusion of the image $f(G)$ as a subgroup of H .

Theorem 4.50 (Second isomorphism theorem). If H is a subgroup of G and N is a normal subgroup of G , then

$$H/(H \cap N) \cong (HN)/N.$$

Proof. Let $f : H \rightarrow (HN)/N$ be defined by $f(x) = xN$. This f is a group homomorphism since $f(xy) = xyN = xNyN = f(x)f(y)$.

Next, we'll show that f is an epimorphism. Let $xN \in (HN)/N$ where $x \in HN$. Then $x = yz$ for some $y \in H$ and $z \in N$. So $xN = yzN = yN = f(y)$. Thus, f is an epimorphism, that is, $f(H) = (HN)/N$. by the first isomorphism theorem, we have

$$H/\text{Ker } f \cong (HN)/N.$$

Finally, we'll show that $\text{Ker } f = H \cap N$ which will imply $H/(H \cap N) \cong (HN)/N$. Let x be an element of H which lies in $\text{Ker } f$. Then xN is the identity element N in $(HN)/N$, so $x \in N$. But $x \in H$ also, so $x \in H \cap N$. Conversely, $x \in H \cap N$ implies $x \in \text{Ker } f$. Q.E.D.

Theorem 4.51 (Third isomorphism theorem). If H and K are both normal subgroups of G with $H \subseteq K$, then

$$(G/H)/(K/H) \cong G/K.$$

Exercise 74. Prove the third isomorphism theorem. Define $f : G/H \rightarrow G/K$ by $f(aH) = aK$. Check that this is a well-defined homomorphism. Show $\text{Ker } f = H$. Show the image of f is all of G/K . Apply the first isomorphism theorem to finish the proof.

Theorem 4.52 (Correspondence theorem). Let N be a normal subgroup of G . The subgroups of G containing N are in one-to-one correspondence with the subgroups of G/N . Thus, if H is a subgroup of G containing N , then H/N is a subgroup of G/N , and every subgroup of G/N so arises. Furthermore, H is normal in G if and only if H/N is normal in G/N .

Exercise 75. Prove the correspondence theorem. Show that for $H \supseteq N$ that H/N is, indeed, a subgroup of G/N . Show that if \overline{H} is any subgroup of G/N that the set $H = \{x \in G \mid x/N \in \overline{H}\}$ is a subgroup of G containing N . Verify that these two operations are inverse to each other. Finally, verify the last statement.

4.6.4 Internal direct products

We can recognize when a group G is isomorphic to a product of two or more groups. Recall that if $G = M \times N$, then we can interpret M and N as subgroups of G . As such they are normal subgroups of G and their intersection is trivial. Furthermore, $G = MN$.

Definition 4.53. A group G is said to be an *internal direct product* of two subgroups M and N if $M \cap N = 1$, $MN = G$, and both M and N are normal subgroups of G .

We'll show in a moment that if G is the internal direct product of M and N , then G is isomorphic to the product group $M \times N$. But first, a lemma.

Lemma 4.54. If M and N are two normal subgroups of G whose intersection is trivial, then elements of M commute with elements of N .

Proof. Let $m \in M$ and $n \in N$. In order to show that $mn = nm$, we'll show the equivalent $mnm^{-1}n^{-1} = 1$. Let $x = mnm^{-1}n^{-1}$. Since $x = (mnm^{-1})n^{-1}$, and both mnm^{-1} and n^{-1} are elements of the normal subgroup N , therefore $x \in N$. But since $x = m(nm^{-1}n^{-1})$, and both m and $nm^{-1}n^{-1}$ are elements of the normal subgroup M , therefore $x \in M$. Since $x \in M \cap N = 1$, therefore $x = 1$. Q.E.D.

Theorem 4.55. If G is the internal direct product of M and N , then $M \times N \cong G$ where the isomorphism is given by $(m, n) \mapsto mn$.

Proof. Outline. Use the lemma to verify that the proposed isomorphism is a homomorphism. It's evidently a surjection since $MN = G$. To show that it's an injection, show that the kernel is trivial. Suppose $(m, n) \mapsto mn = 1$. Then $m = n^{-1}$ lies in both M and N , so it's trivial, that is, $m = n = 1$. Q.E.D.

Exercise 76. Prove that G is an internal direct product of two normal subgroups M and N if and only if every element $x \in G$ can be uniquely represented as a product mn with $m \in M$ and $n \in N$.

Although we've only looked at internal direct products of two subgroups, the definition can be generalized to more than two subgroups. We'll say that G is the *internal direct product* of r normal subgroups N_1, N_2, \dots, N_r if (1) they jointly generate G , that is, $N_1N_2 \cdots N_r = G$, and (2) the intersection of any one N_i with the subgroup generated by the rest is trivial. It follows that $N_1 \times N_2 \times \cdots \times N_r \cong G$. Furthermore, an equivalent condition to being an internal direct product of the normal subgroups N_1, N_2, \dots, N_r is that every element $x \in G$ can be uniquely represented as a product $n_1n_2 \cdots n_r$ with each $n_i \in N_i$.

4.7 Matrix rings and linear groups

The representation of rings and groups as subrings or subgroups of matrix rings is very helpful for a couple of reasons. One is that matrices describe linear transformations. That means that the elements of the ring or group can be interpreted as geometric transformations. A second is that matrix notation is so very convenient. Usually the coefficients are taken to be elements of a familiar field like \mathbf{C} , \mathbf{R} , or \mathbf{Q} , but for special purposes the coefficients may be taken in some other integral domain such as \mathbf{Z} .

For example, the field complex numbers \mathbf{C} can be represented as a certain subring of $M_2(\mathbf{R})$, the ring of 2×2 matrices with coefficients in \mathbf{R} , and the division ring of quaternions \mathbf{H} can be represented as a certain subring of $M_4(\mathbf{R})$.

Most of our examples have n equal to 2 or 3 and the coefficients are real.

4.7.1 Linear transformations

The ring of $n \times n$ matrices with real coefficients, $M_2(\mathbf{R})$, is a noncommutative ring when $n \geq 2$. We can interpret each matrix $A \in M_2(\mathbf{R})$ as a linear transformation $A : \mathbf{R}^n \rightarrow \mathbf{R}^n$

where a (column) n -vector $\mathbf{x} \in \mathbf{R}^n$ is mapped to another n -vector

$$A\mathbf{x} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n \\ \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n \end{bmatrix}$$

The identity matrix

$$I = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}$$

corresponds to the identity transformation $I : \mathbf{R}^n \rightarrow \mathbf{R}^n$ where $I\mathbf{x} = \mathbf{x}$.

A linear transformation from a vector space to itself is also called a linear operator.

4.7.2 The general linear groups $GL_n(R)$

The invertible $n \times n$ matrices in $M_n(R)$, that is, the units in the ring $M_n(R)$, form the *general linear group* with coefficients in the commutative ring R , denoted $GL_n(R)$. They describe nonsingular transformations $R^n \rightarrow R^n$. Recall that a matrix A has an inverse if and only if its determinant $|A|$ is a unit in R .

Let's interpret some of these in the case when $n = 2$. The determinant of $A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ is $|A| = ad - bc$, and when that's a unit in R , the inverse of A is $A = \frac{1}{|A|} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$.

Note that the determinant is a group homomorphism $GL_n(R) \rightarrow R^*$ from the general linear group to the invertible elements of R . The determinant of the identity matrix is 1, the determinant of the product of two matrices is the product of their determinants, and the determinant of the inverse of a matrix is the reciprocal of the determinant of the matrix.

Let's let R be the field of real numbers \mathbf{R} . The real general linear group $GL_2(\mathbf{R})$ can be interpreted as the group of invertible linear transformations of the plane \mathbf{R}^2 that leave the origin fixed. Here are a few linear transformations of the plane.

Rotation by an angle θ about the origin is described by the matrix

$$\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}$$

since a point $\begin{bmatrix} x \\ y \end{bmatrix}$ in \mathbf{R}^2 is sent to the point $\begin{bmatrix} x \cos \theta - y \sin \theta \\ x \sin \theta + y \cos \theta \end{bmatrix}$. The determinant of a rotation matrix is 1.

Reflection across a line through the origin at an angle θ to the x -axis is described by the matrix

$$\begin{bmatrix} \cos 2\theta & \sin 2\theta \\ \sin 2\theta & -\cos 2\theta \end{bmatrix}.$$

The determinant is -1 .

Expansions and contractions are described by scalar matrices $\begin{bmatrix} r & 0 \\ 0 & r \end{bmatrix}$ where r is the ratio. If $r > 1$, then it's an expansion (also called dilation), but if $0 < r < 1$, then it's a contraction.

There are numerous other kinds of transformations. Here's just one more example $\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$, an example of a shear parallel to the x -axis. Points above the x -axis are moved right, points below left, and points on the x -axis are fixed.

In three dimensions you can describe rotations, reflections, and so forth, as well.

4.7.3 Other linear groups

There are a number of interesting subgroups of $GL_n(R)$.

The special linear groups $SL_n(R)$. There are several subgroups of $GL_n(R)$, one of which is the special linear group $SL_n(R)$ which consists of matrices whose determinants equal 1, also called *unimodular* matrices. (There are other linear groups called "special" and in each case it means the determinant is 1.)

Among the examples in $GL_2(\mathbf{R})$ mentioned above, the rotations and shears are members of $SL_2(\mathbf{R})$, but reflections have determinant -1 and expansions and contractions have determinants greater or less than 1, so none of them belong to the special linear group.

Since the absolute value of the determinant is the Jacobian of the transformation $\mathbf{R}^n \rightarrow \mathbf{R}^n$, therefore transformations in $SL_2(\mathbf{R})$ preserve area. Since the determinant is positive, these transformations preserve orientation. Thus, transformations in $SL_2(\mathbf{R})$ are the linear transformations that preserve orientation and area. More generally those in $SL_n(\mathbf{R})$ preserve orientation and n -dimensional content. Rotations and shears, and their products, are always in $SL_n(\mathbf{R})$.

Exercise 77. Show that the matrix $\begin{bmatrix} 2 & 0 \\ 0 & 1/2 \end{bmatrix}$ lies in $SL_2(\mathbf{R})$. Describe in words how this transformation acts on the plane.

The orthogonal groups $\mathcal{O}(n)$. These are subgroups of $GL_n(\mathbf{R})$. An *orthogonal* transformation is one that preserves inner products (also called dot products or scalar products).

I'll use the notation

$$\langle \mathbf{a} | \mathbf{b} \rangle = a_1 b_1 + a_2 b_2 + \cdots + a_n b_n$$

for the inner product of the vectors $\mathbf{a} = (a_1, a_2, \dots, a_n)$ and $\mathbf{b} = (b_1, b_2, \dots, b_n)$. Other common notations are (\mathbf{a}, \mathbf{b}) or $\mathbf{a} \cdot \mathbf{b}$.

For the transformation described by the matrix A to preserve inner products means that $\langle A\mathbf{a} | A\mathbf{b} \rangle = \langle \mathbf{a} | \mathbf{b} \rangle$. Since the length of a vector $|\mathbf{a}|$ is determined by the inner product, $|\mathbf{a}|^2 = \langle \mathbf{a} | \mathbf{a} \rangle$, therefore an orthogonal transformation preserves distance, too: $|A\mathbf{a}| = |\mathbf{a}|$. Conversely, if A preserves distance, it preserves inner products.

Note that since distance is preserved, so is area in dimension 2 or n -dimensional content in dimension n .

It's a theorem from linear algebra that a matrix A describes an orthogonal transformation if and only if its inverse equals its transform: $A^{-1} = A^T$; equivalently, $AA^T = 1$. These ma-

trices, of course, are called *orthogonal* matrices. Note that the determinant of an orthogonal matrix is ± 1 .

The orthogonal group $\mathcal{O}(n)$ is the subgroup of $GL_n(\mathbf{R})$ of orthogonal matrices. It's not a subgroup of $SL_n(\mathbf{R})$ since half the orthogonal matrices have determinant -1 , meaning they reverse orientation. The special orthogonal group $S\mathcal{O}(n)$ is the subgroup of $\mathcal{O}(n)$ of matrices with determinant 1.

In two dimensions $\mathcal{O}(2)$ consists of rotations and reflections while $S\mathcal{O}(2)$ consists of only the rotations. In three dimensions $\mathcal{O}(3)$ consists of rotations (by some angle around some line through 0) and reflections (across some plane through 0). Again, $S\mathcal{O}(3)$ only has the rotations.

The unitary groups $\mathcal{U}(n)$. For matrices with complex coefficients, the most useful analogous group corresponding to the orthogonal group for real coefficients is something called a unitary group.

The inner product, also called the *Hermitian*, for the complex vector space \mathbf{C}^n is defined as

$$\langle \mathbf{a} | \mathbf{b} \rangle = a_1 \bar{b}_1 + a_2 \bar{b}_2 + \cdots + a_n \bar{b}_n$$

for the complex vectors $\mathbf{a} = (a_1, a_2, \dots, a_n)$ and $\mathbf{b} = (b_1, b_2, \dots, b_n)$ where the bar indicates complex conjugation. A matrix A , and the transformation $\mathbf{C}^n \rightarrow \mathbf{C}^n$ that it describes, are called *unitary* if it preserves the Hermitian. The collection of all unitary matrices in $GL_n(\mathbf{C})$ is called the unitary group $\mathcal{U}(n)$.

Another theorem from linear algebra is that a matrix A is unitary if and only if its inverse is the transform of its conjugate, $A^{-1} = \bar{A}^T$, equivalently, $A\bar{A}^T = I$.

There are many properties of complex unitary matrices that correspond to properties of real orthogonal matrices.

4.7.4 Projective space and the projective linear group $PGL_n(F)$

Projective planes and projective space. Projective geometry differs from Euclidean geometry in a couple of ways: all lines in a plane intersect, and distance and angles are not considered.

Let's start with Euclidean plane geometry, then drop distance and angles, then add points at infinity to get the projective plane.

When distance and angles are not considered in Euclidean geometry, what's left is called *affine geometry*. Points and lines still remain. The affine plane is still modelled by \mathbf{R}^2 , but affine transformations don't have to preserve distance or angles. So, for instance, the linear transformations known as expansions, contractions, and shear transformations are all affine transformations. In fact every element in $GL_2(\mathbf{R})$ describes an affine planar transformation. These are the affine transformations that fix the origin. Also, translations, which are not linear transformations, are affine transformations. Similarly, in dimension n , affine transformations are composed of translations and elements of $GL_n(\mathbf{R})$.

Affine spaces F^n can be similarly defined for other fields F besides the reals \mathbf{R} .

So far, we've dropped distance and angles, but parallel lines remain in affine geometry. The next step is to add enough points, called points at infinity, so that parallel lines meet at them.

Parallelism is an equivalence relation on lines. We'll assume that a line is parallel to itself, so parallelism is reflexive. It's also symmetric, and it's transitive: if one line is parallel to another, and the other to a third, then the first is parallel to the third.

For each parallelism equivalence class, add one point, a *point at infinity* to affine space and specify that every line in that equivalence class passes through that point. Add one more line, the *line at infinity*, and specify that every point at infinity passes through it. The resulting space is the *projective space* corresponding to the affine space.

Projective space and projective coordinates. Let F be a field, such as the field of real numbers. The projective linear group $PGL_n(F)$ is used to study projective space.

A more formal way to define projective n -space over a field F is by modelling points of the projective plane by lines in affine $n + 1$ -space, F^{n+1} , through the origin by means an algebraic equivalence relation.

Two points $\mathbf{a} = (a_0, a_1, \dots, a_n)$ and $\mathbf{b} = (b_0, b_1, \dots, b_n)$ of F^{n+1} name the same point of FP^n if their coordinates are proportional, that is, if there exists a nonzero element $\lambda \in F$ such that $b_i/a_i = \lambda$ for $i = 0, 1, \dots, n$. We'll let $[a_0, a_1, \dots, a_n]$ denote the point in FP^n named by $(a_0, a_1, \dots, a_n) \in F^{n+1}$. Thus, $[a_0, a_1, \dots, a_n] = [\lambda a_0, \lambda a_1, \dots, \lambda a_n]$. The notation $[a_0, a_1, \dots, a_n]$ is called *projective coordinates*.

Geometrically, this construction adds points at infinity to the affine plane, one point for each set of parallel lines.

Lines can also be named with projective coordinates $\mathbf{b} = [b_0, b_1, \dots, b_n]$. If you do that, then a point $\mathbf{a} = [a_0, a_1, \dots, a_n]$ lies on the line \mathbf{b} if their inner product $\langle \mathbf{a} | \mathbf{b} \rangle$ is 0.

Example 4.56 (The Fano plane \mathbf{Z}_2P^2). The projective plane \mathbf{Z}_2P^2 has a name, the Fano plane, named after Gino Fano (1871–1952), a founder of finite geometries.

Figure 4.9 shows a representation of \mathbf{Z}_2P^2 . There are 7 points and 7 lines, each line with 3 points, and each point on 3 lines.

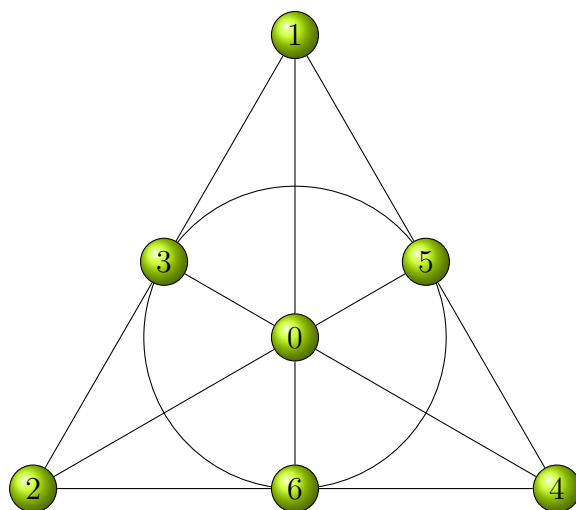


Figure 4.9: The Fano plane \mathbf{Z}_2P^2

Example 4.57 (The projective plane \mathbf{Z}_3P^2). Figure 4.10 shows a representation of the finite projective plane \mathbf{Z}_3P^2 . There are 13 points and 13 lines, each line with 4 points, and each point on 4 lines.

We can name the 9 points in the affine plane \mathbf{Z}_3^2 with third coordinate 1, and the 4 points at infinity with third coordinate 0. The four points at infinity lie on a line at infinity. Each of these points at infinity lie on all those line with a particular slope. For instance, the point $[1, -1, 0]$ lies on the three lines with slope -1 (and it lies on the line at infinity, too).

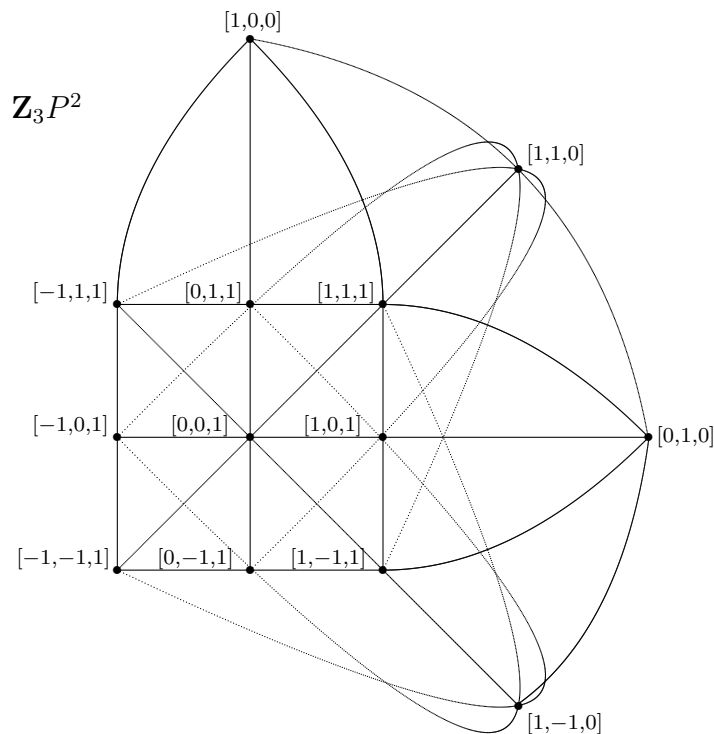


Figure 4.10: The projective plane \mathbf{Z}_3P^2

Finite projective planes. There's a simpler definition of a projective plane that can be made axiomatically. It states that two points determine a line, and two lines determine a point. A nondegeneracy axiom is also required that there are at least three points which don't all lie on the same line (from which it follows that there are at least three lines which don't all meet at one point). It turns out that this axiomatic definition admits projective planes that don't derive from fields. We'll look at the ones that do.

Finite projective planes exist for each finite field. Let $GF(p^n)$ be a Galois field of $q = p^n$ elements. There will be q^2 points on the affine plane $GF(p^n)^2$ with third coordinate 1, and $q + 1$ points on the line at infinity with third coordinate 0. So the finite projective plane $GF(p^n)P^2$ has $q^2 + q + 1$ points altogether. It has the same number of lines.

These projective planes all have a couple of nice properties. They are all Desarguesian and Pappian, that is, Desargue's theorem and Pappas's theorem both hold for these projective planes. These two theorems state that certain configurations of points and lines hold for the

projective plane. Desargues developed projective geometry in the 1600s, and one of Pappus's theorems apply to projective geometry. There are other projective planes that aren't based on finite fields that aren't Desarguesian and Pappian.

Projective linear group $PGL_n(F)$. As we defined projective $n - 1$ -space over a field F as a quotient of nonzero elements of n -space, so too we can define a quotient of $GL_n(F)$ to get the projective linear group $PGL_n(F)$ acting on projective $n - 1$ -space. Two matrices A and B in $GL_n(F)$ name the same element of $PGL_n(F)$ if each is a multiple of the other, that is, there exists $\lambda \neq 0 \in F$ such that $B = \lambda A$. Then $PGL_n(F)$ acts on FP^{n-1} , since $A\mathbf{a}$ and $\lambda A\mathbf{a}$ name the same element of FP^{n-1} .

If F is a finite field with q elements, then the order of the group $PGL_n(F)$ is the order of $GL(n, F)$ divided by $q - 1$, so $|PGL_n(F)| = \frac{(q^n - 1)(q^n - q)(q^n - q^2) \cdots (q^n - q^{n-1})}{q - 1}$.

Projective special linear group $PSL_n(F)$. The *projective special linear group*, $PSL_n(F)$, is the subgroup of $PGL_n(F)$ named by unimodular matrices. It's $SL_n(F)$ modulo scalar matrices ωI where ω is an n th root of unity.

The order of $PSL_n(F)$ is equal to the order of $PGL_n(F)$ divided by $\text{GCD}(n, q - 1)$ where q is the number of elements of F .

Except for small values of n the projective special linear groups are all simple. Simplicity is defined in the next section.

The groups $PSL_3(\mathbf{Z}_3)$ is actually the same as $PGL_3(\mathbf{Z}_3)$ since 3 and 2 are relatively prime.

Example 4.58. The projective linear group $PGL_3(\mathbf{Z}_2) = PSL_3(\mathbf{Z}_2)$ which acts on the Fano plane \mathbf{Z}_2P^2 has $\frac{7 \cdot 6 \cdot 4}{1} = 168$ elements.

It's small enough so that its conjugacy classes can be determined without resorting to advanced methods. There are six conjugacy classes of sizes 1, 21, 56, 42, 24, and 24. As always, the identity forms a conjugacy class of 1 element. Refer to figure 4.9 to name elements. The conjugacy class of the involution (12)(56) has size 21. The conjugacy class of (124)(365) has 56 elements. The conjugacy class of (0124)(36) has 42 elements. The conjugacy class of (0125463) has 24 elements, and the conjugacy class of its inverse also has 24 elements.

Also, $PGL_3(\mathbf{Z}_3) = PSL_3(\mathbf{Z}_3)$, acting on the projective plane \mathbf{Z}_3P^2 , has order $\frac{26 \cdot 24 \cdot 18}{2} = 5616$.

4.8 Structure of finite groups

The classification of finite groups is extremely difficult, but there are a tools we can use to see how that classification begins. In the next section we'll classify finite Abelian groups and see that they're isomorphic to products of cyclic groups, but the situation for general groups much more complicated.

4.8.1 Simple groups

The way we'll analyze groups is by their normal subgroups and quotients. In particular, if N is a maximal, proper normal subgroup of G , then G/N has no subgroups, for if it did, by the correspondence theorem, there would be a normal subgroup between N and G .

Definition 4.59. A nontrivial group is said to be *simple* if it has no proper, nontrivial, normal subgroups.

Exercise 78. Prove that the only Abelian simple groups are cyclic of prime order.

There are many nonabelian simple groups. There are several infinite families of them, and a few that aren't in infinite families, called *sporadic* simple groups. One infinite family of simple groups consists of alternating groups A_n with $n \geq 5$. Indeed, A_5 is the smallest nonabelian simple group. The projective special linear groups mentioned in the section above form another family of finite simple groups.

Exercise 79 (Nonsimplicity of A_4). Verify that there are five conjugacy classes in A_4 as shown in the following table.

Generator	Size	Order
1	1	1
(12)(34)	3	2
(123)	4	3
(132)	4	3

A normal subgroup of A_4 would be a union of some of these conjugacy classes including the identity conjugacy class of size 1, but its order would have to divide 12. Find all the proper nontrivial normal subgroups of A_4 .

Exercise 80 (Simplicity of A_5). Verify that there are five conjugacy classes in A_5 as shown in the following table.

Generator	Size	Order
1	1	1
(12)(34)	15	2
(123)	20	3
(12345)	12	5
(12354)	12	5

A normal subgroup of A_5 would be a union of some of these conjugacy classes including the identity conjugacy class of size 1, but its order would have to divide 60. Verify that no combination of the numbers 1, 15, 12, 12, and 20, where 1 is included in the the combination, yields a sum that divides 60 (those numbers being 2, 3, 4, 6, 10, 12, 15, 20, and 30) except just 1 itself and the sum of all five numbers. Thus, there is no proper nontrivial normal subgroup of A_5 .

4.8.2 The Jordan-Hölder theorem

Definition 4.60. A *composition series* for a group G is a finite chain of subgroups

$$1 = N_n \subseteq N_{n-1} \subseteq \cdots \subseteq N_1 \subseteq N_0 = G$$

such that each N_{i-1} is a maximal proper normal subgroup of N_i . The number n is called the *length* of the composition series, and the n quotient groups

$$N_{n-1}/1, \dots, N_1/N_2, G/N_1$$

which are all simple groups, are called *composition factors* determined by the composition series.

It is evident that any finite group G has at least one composition series. Just take N_1 to be a maximal proper normal subgroup of G , N_2 to be a maximal proper normal subgroup of N_1 , etc. Infinite groups may also have composition series, but not all infinite groups do.

Exercise 81. Find a composition series for the symmetric group S_4 .

Exercise 82. Prove that an infinite cyclic group has no (finite) composition series.

Although a finite group may have more than one composition series, the length of the series is determined by the group as are composition factors at least up to isomorphism as we'll see in a moment. Thus, these are invariants of the group. They do not, however, completely determine the group.

Exercise 83. Show that the dihedral group D_5 and the cyclic group C_{10} have composition series with the same length and same factors.

Theorem 4.61 (Jordan-Hölder). Any two composition series for a finite group have the same length and there is a one-to-one correspondence between the composition factors of the two composition series for which the corresponding composition factors are isomorphic.

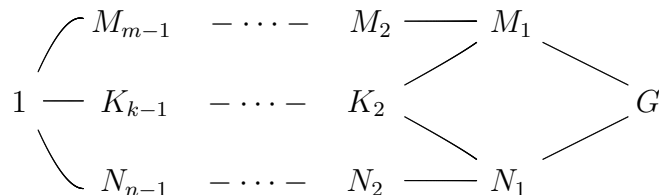
Proof. We'll prove this by induction on the order of the group under question. The base case is for the trivial group which has only the trivial composition series.

Assume now that a group G has two composition series

$$1 = N_m \subseteq M_{m-1} \subseteq \dots \subseteq M_1 \subseteq M_0 = G, \text{ and } 1 = N_n \subseteq N_{n-1} \subseteq \dots \subseteq N_1 \subseteq N_0 = G$$

If $M_1 = N_1$, then by induction we conclude that the lengths of the rest of the composition are equal and the composition factors the rest of the rest of the series are the same, and of course, the factors G/M_1 and G/N_1 are equal, so the case $M_1 = N_1$ is finished.

Consider now the case $M_1 \neq N_1$. Since both M_1 and N_1 are normal subgroups of G , so is their intersection $K_2 = M_1 \cap N_1$. Let $1 = K_k \subseteq K_{k-1} \subseteq \dots \subseteq K_3 \subseteq K_2$ be a composition series for their intersection. These subgroups of G are illustrated in the following diagram.



By the second isomorphism theorem, we have $M_1/(M_1 \cap N_1) \cong G/N_1$. Therefore, K_2 is a maximal normal subgroup of M_1 . Thus, we have two composition series for M_1 , and by the inductive hypothesis, they have the same length, so $m = k$, and they have the same factors

up to isomorphism in some order. Likewise we have two composition series for N_1 , and they have the same length, so $k = n$, and the same factors up to isomorphism in some order. We now have four composition series for G , two including M_1 and two including N_1 . They all have the same length, and since $G/M_1 \cong N_1/K_2$ and $G/N_1 \cong M_1/K_2$, they all have the same factors up to isomorphism in some order. Q.E.D.

There is a generalization of this theorem that applies to infinite groups that have composition series but its proof is considerably longer.

The list of composition factors is not enough to characterize the group. That is to say, there are non-isomorphic groups that have the same composition factors. The smallest pair of such groups are A_3 and C_6 of order 6.

A sporadic group. Most finite simple groups come in infinite parameterized families such as the cyclic groups C_p for prime p , and the alternating groups A_n for $n \geq 5$. There are several of these infinite families of simple groups. There are also a few simple groups that don't belong to any of these infinite families. We'll look at one of them, the Mathieu group M_{11} .

Mathieu discovered M_{11} in 1861. It's the smallest sporadic group, and it has order $7920 = 8 \cdot 9 \cdot 10 \cdot 11$. It can be described as a subgroup of the symmetric group S_{11} generated by the pair of permutations $(123456789te)$ and $(37e8)(4t56)$. (Here t is used for 10 and e for 11.)

M_{11} has elements of order 1, 2, 3, 4, 5, 6, 8, and 11. It has $165 = 3 \cdot 5 \cdot 11$ elements of order 2, that is, involutions. They are all conjugates of $(3t)(49)(56)(8e)$.

As S_{11} acts on a set of 11 elements, so does M_{11} . In fact, the action is sharply 4-transitive. *Transitive* means that for any pair x and y of elements in the set, there is a group element that maps the x to y . *Doubly transitive* means that for x_1, x_2 and y_1, y_2 , distinct pairs, there's a group element that sends x_1 to y_1 at the same time as sending x_2 to y_2 . More generally, and *n-transitive action* is one such that for all pairwise distinct n -tuples x_1, \dots, x_n and pairwise distinct y_1, \dots, y_n there is a group element that maps each x_i to y_i . When there is exactly one group element for pair of n -tuples, the group is said to act *sharply*.

Solvable groups One of the applications of group theory is Galois theory for algebraic fields. The groups of automorphisms of these fields are closely related to the solutions of algebraic equations. In particular, these groups can tell you if the equations have solutions that can be expressed in terms of radicals, that is square roots, cube roots, and higher roots. The condition for such solvability is none the factors in a composition series for a group are nonabelian simple groups, equivalently, that all the factors are cyclic groups of prime order.

Definition 4.62. A group is said to be *solvable* if it has a composition series all of whose factors are cyclic.

Exercise 84. Prove that if the order of a group is a power of a prime number, then that group is solvable.

Example 4.63 (The Frobenius group $F_{21} = C_7 \rtimes C_3$). This group will have 21 elements. It is what is called a semidirect product of the cyclic group $C_7 = \{1, a, a^2, \dots, a^6\}$ of 7 elements with the cyclic group $C_3 = \{1, b, b^2\}$ of 3 elements. Each element can be written in the form $b^n a^m$ with $0 \leq b \leq 2$ and $0 \leq a \leq 6$, but a and b don't commute. For this group, $bab^{-1} = a^2$.

The group is denoted $C_7 \rtimes C_3$. The group C_7 is a normal subgroup, but C_3 is not a normal subgroup.

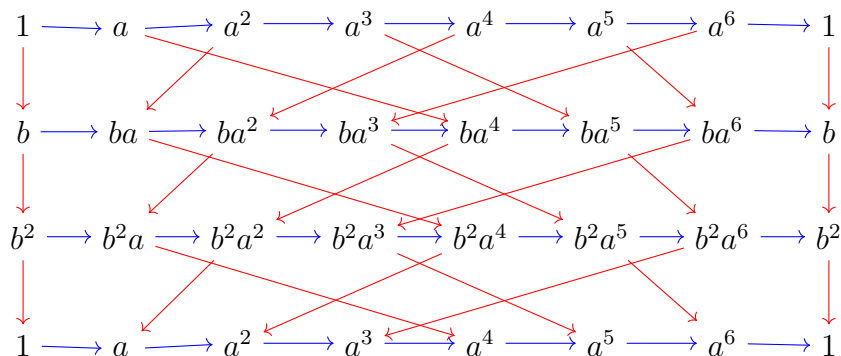


Figure 4.11: Cayley graph of the Frobenius group $F_{21} = C_7 \rtimes C_3$

This group can be presented as $\langle a, b : a^7 = 1, b^3 = 1, ba = a^2b \rangle$.

Its Cayley graph is shown in figure 4.11 with the understanding that the top line is identified with the bottom line, and the left line is identified with the right line. A blue arrow indicates multiplication by a while a red one is multiplication by b .

The group $C_7 \rtimes C_3$ is a group of symmetries of a heptahedron on a torus.

A heptahedron has 7 hexagonal faces which meet three at a time at a vertex, 14 vertices, and 27 edges. It is a tiling of the torus which is illustrated in figure 4.12. Each of the seven hexagons is labelled 1 through 7 and colored a different color. The outer edges are to be identified so that the edges $ABCD$ are identified on the upper left and lower right, the edges $DEFA$ are identified on the upper left and lower right, and the edges $AGHD$ are identified on the left and the right. The resulting topological space is a torus.

You can also interpret this as a coloring of the tiling of the plane by hexagons where the labels of some nearby hexagons are shown in the figure.

The group $C_7 \rtimes C_3$ is a subgroup of the group of symmetries of this heptahedron. The element a of order 7 describes the permutation of the faces (1234567) which moves the hexagons to the upper right. The element b of order 3 describes the permutation (142)(356) which is a rotation about hexagon 7 by 120° . It's easily verified that ba and a^2b both describe the permutation (157)(364) which is a rotation about hexagon 2.

Exercise 85. Verify that the rotation $c = (154623)$ about hexagon 7 by 60° is a symmetry of the heptahedron. Evidently $c^2 = b$.

- Determine the relation between a and c of the form $ca = a^k c$, that is, find k .
- This group is a semidirect product $C_7 \rtimes C_6$. Draw its Cayley graph.

Much more can be said about solvable groups than we have time for.

4.9 Abelian groups

Commutative groups are called Abelian groups in honor of Neils Henrik Abel (1802–1829) who worked with groups of substitutions in order to understand solutions of polynomial equations.

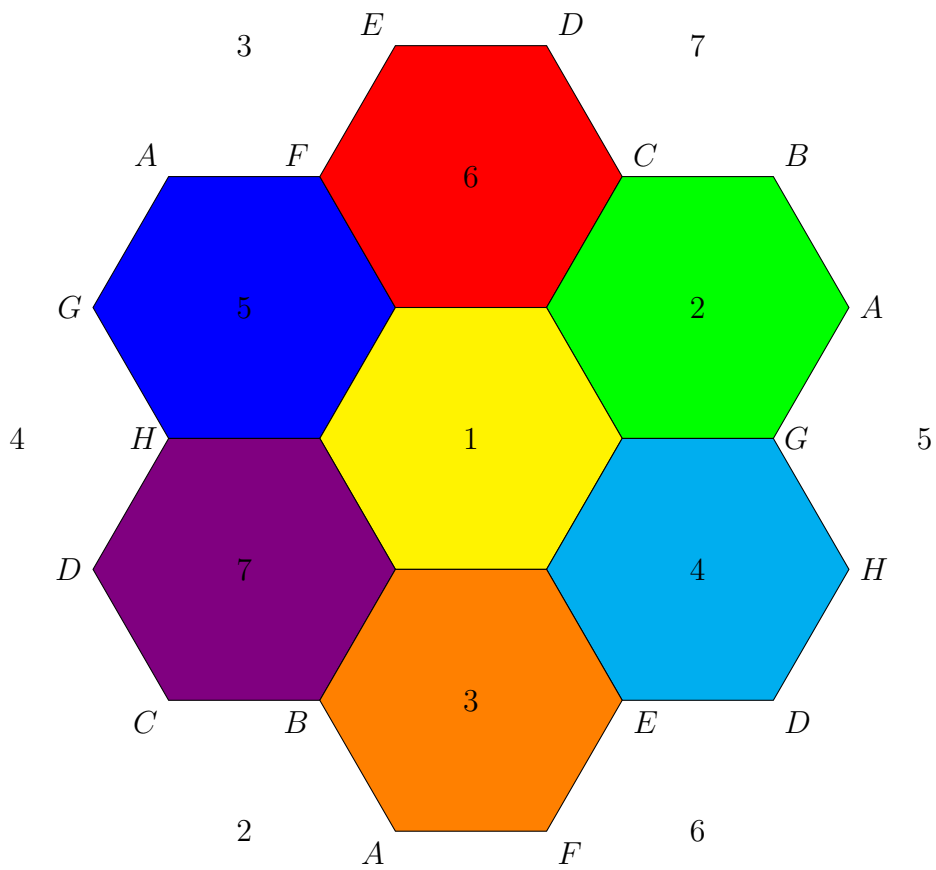


Figure 4.12: Heptahedron on a torus

We'll use additive notation throughout this section on Abelian groups. Also, we'll call the product of two Abelian groups A and B a direct sum and denote it $A \oplus B$ rather than $A \times B$.

Every subgroup of an Abelian group is normal, so we'll just refer to them as subgroups and leave off the adjective "normal."

We already know a fair amount about Abelian groups. We know about cyclic groups and the Chinese remainder theorem.

For example, we know $\mathbf{Z}_{12} \cong \mathbf{Z}_3 \oplus \mathbf{Z}_4$ where an element n modulo 12 corresponds to the pair n modulo 3 and n modulo 4. Likewise, $\mathbf{Z}_6 \cong \mathbf{Z}_2 \oplus \mathbf{Z}_3$. This gives us three ways to treat the group $\mathbf{Z}_2 \oplus \mathbf{Z}_3 \oplus \mathbf{Z}_4$ since it is isomorphic to both $\mathbf{Z}_2 \oplus \mathbf{Z}_{12}$ and $\mathbf{Z}_6 \oplus \mathbf{Z}_4$.

Our characterization of internal direct product looks a little different when the group is written additively. Here it is, rewritten for Abelian groups.

An Abelian group G is the *internal direct sum* of subgroups M and N if (1) they jointly generate G , that is, $M + N = G$, and (2) the intersection $M \cap N = 0$. If G is the internal direct sum of M and N , then $M \oplus N = G$. Furthermore, an equivalent condition to being an internal direct sum is that every element $x \in G$ can be uniquely represented as a sum $m + n$ with $m \in M$ and $n \in N$.

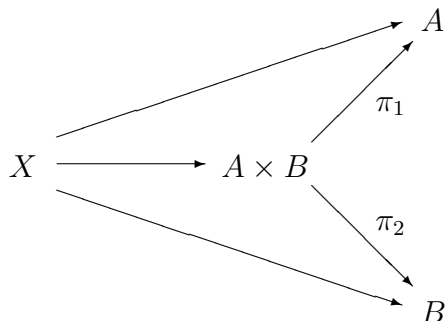
For the example $\mathbf{Z}_2 \oplus \mathbf{Z}_3 \oplus \mathbf{Z}_4$ above, it is the internal direct sum of \mathbf{Z}_2 and $0 \oplus \mathbf{Z}_3 \oplus \mathbf{Z}_4$ as well as the internal direct sum of $\mathbf{Z}_2 \oplus \mathbf{Z}_3 \oplus 0$ and \mathbf{Z}_4 .

4.9.1 The category \mathcal{A} of Abelian groups

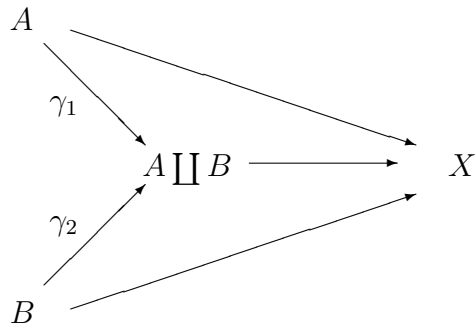
The category of Abelian groups is a particularly nice category. Not only does it have products, but it also has coproducts, to be defined next, and the products are coproducts, and that's why we're calling them direct sums. It's not the only category with direct sums. The category of vector spaces over a fixed field has them too.

Coproducts in a category and their universal property When all the arrows in a diagram are reversed, a similar diagram, called the *dual* results. Recall that products in a category are characterized by a diagram.

The product $A \times B$ in a category along with the two projections $A \times B \xrightarrow{\pi_1} A$ and $A \times B \xrightarrow{\pi_2} B$ has the universal property that for each object X and morphisms $X \rightarrow A$ and $X \rightarrow B$, there is a unique morphism $X \rightarrow A \times B$, such that the diagram below commutes.



If we turn around all the arrows, we'll get the characterizing property for coproducts. The coproduct $A \coprod B$ in a category along with the two injections $A \xrightarrow{\gamma_1} A \coprod B$ and $B \xrightarrow{\gamma_2} A \coprod B$ has the universal property that for each object X and morphisms $A \rightarrow X$ and $B \rightarrow X$, there is a unique morphism $A \coprod B \rightarrow X$, such that the diagram below commutes.



In the category \mathcal{S} of sets coproducts are disjoint unions. The disjoint union of two sets S and T has one element for each element of S and a different element for each element of T . So the cardinality of their disjoint union is $|S| + |T|$.

Exercise 86. In the category of Abelian groups, the coproduct object $A \coprod B$ is what we've called the direct sum $A \oplus B$, which is the same as the product $A \times B$. The injections $A \xrightarrow{\gamma_1} A \coprod B$ and $B \xrightarrow{\gamma_2} A \coprod B$ for Abelian groups are defined by $\gamma_1(x) = (x, 0)$ and $\gamma_2(y) = (0, y)$. Verify that the universal property holds.

4.9.2 Finite Abelian groups

The classification of finite groups is very difficult, but the classification of finite Abelian is not so difficult. It turns out, as we'll see, that a finite Abelian group is isomorphic to a product of cyclic groups, and there's a certain uniqueness to this representation. This classification is sometimes called the fundamental theorem of finite Abelian groups. The theorem above on internal direct sums is essential in this classification.

Theorem 4.64. Let G be a finite Abelian group of order mn where m and n are relatively prime, both greater than 1. Let $M = \{x \in G \mid mx = 0\}$ and $N = \{x \in G \mid nx = 0\}$. Then M and N are subgroups of G , and G is the internal direct sum of M and N . Furthermore, $|M| = m$ and $|N| = n$.

Proof. Outline. That M and N are subgroups is quickly verified. Since m and n are relatively prime, therefore 1 is a linear combination of them, that is, there are integers s and t such that $1 = sm + tn$. Their intersection $M \cap N$ is trivial since if $x \in M \cap N$, then $mx = nx = 0$, hence $x = 1x = (sm + tn)x = smx + tnx = 0$. Together M and N generate G , since for $x \in G$, $x = smx + tnx$, but $smx \in N$ since $nsmx = (nm)sx = 0$, likewise $tnx \in M$. Thus $M + N = G$. Therefore, G is the internal direct sum of M and N . Q.E.D.

p -primary groups. Let G be an Abelian group and p a prime number. The set

$$G(p) = \{x \mid p^k x = 0 \text{ for some } k \geq 0\}$$

is a subgroup of G . It is called the p -primary component of G .

As a corollary to the above theorem consider the case when $|G|$ is factored as a power of primes.

Corollary 4.65 (Primary decomposition theorem). Let G be a finite Abelian group whose order has prime factorization $p_1^{e_1} p_2^{e_2} \cdots p_r^{e_r}$. Then G is a direct sum of the p_i -primary components

$$G \cong G(p_1) \oplus G(p_2) \oplus \cdots \oplus G(p_r)$$

and $|G(p_i)| = p_i^{e_i}$ for each i .

We've reduced the problem of classifying finite Abelian groups to classifying those whose orders are powers of a prime p . Such groups are called p -primary groups or simply p -groups. If the power is greater than 1, then there are different groups of that order. For example, there are three distinct Abelian groups of order 125, namely, \mathbf{Z}_{125} , $\mathbf{Z}_{25} \oplus \mathbf{Z}_5$ and $\mathbf{Z}_5 \oplus \mathbf{Z}_5 \oplus \mathbf{Z}_5$. The first has an element of order 125, but the other two don't, while the second has an element of order 25, but the third doesn't. Hence, they are not isomorphic.

We'll see soon that every p -group G is isomorphic to a unique direct sum of cyclic p -groups

$$G \cong \mathbf{Z}_p^{e_1} \oplus \mathbf{Z}_{p^2}^{e_2} \oplus \cdots \oplus \mathbf{Z}_{p^r}^{e_r}$$

where the sum $e_1 + 2e_2 + \cdots + re_r$ is equal to n , where $|G| = p^n$.

Example 4.66. We'll find all the 2-groups of order 32 up to isomorphism. Since $32 = 2^5$, we'll need $e_1 + 2e_2 + \cdots + re_r = 5$. Each solution will give us a way of making a sum of positive integers equal to 5. A *partition* of n is a list of positive integers that sum to n . Here's a table which gives all the partitions of 5 and the associated 2-groups.

5	\mathbf{Z}_{32}
1 + 4	$\mathbf{Z}_2 \oplus \mathbf{Z}_{16}$
2 + 3	$\mathbf{Z}_4 \oplus \mathbf{Z}_8$
1 + 1 + 3	$\mathbf{Z}_2^2 \oplus \mathbf{Z}_4$
1 + 2 + 2	$\mathbf{Z}_2 \oplus \mathbf{Z}_4^2$
1 + 1 + 1 + 2	$\mathbf{Z}_2^3 \oplus \mathbf{Z}_4$
1 + 1 + 1 + 1 + 1	\mathbf{Z}_2^5

Exercise 87. Complete this table of the number of partitions of n up through $n = 10$. Work it out yourself.

n	0	1	2	3	4	5	6	7	8	9	10
	1	1	2	3	5	7					

Fundamental theorem of finite Abelian groups

Our strategy for a p -primary group will be to pick off direct summands containing elements of maximal orders, one at a time. That will show that a p -primary group is a direct sum of cyclic groups whose orders are nonincreasing powers of p . We'll then show those powers of p are determined by the p -primary group.

A difficulty in the proof is that there are many choices to be made resulting in different direct sums, but we'll see that the orders of the cyclic subgroups turns out to be the same no matter how we make the choices.

The proof of the theorem is particularly technical, so we'll separate parts of the proof as lemmas.

Lemma 4.67. Let G be a noncyclic p -primary group and a an element of G of maximal order. Then there is an element b in the complement of $\langle a \rangle$ of order p .

Proof. Let c be an element in the complement of $\langle a \rangle$ of smallest order. Since the order of pc is $1/p$ times the order of c , which is a smaller order than the order of c , therefore pc lies in $\langle a \rangle$. So $pc = ka$ for some integer k . Let p^m denote the $\text{ord } a$, the largest order of any element in G . Then $\text{ord}(ka) \leq p^{m-1}$ since $p^{m-1}(ka) = p^{m-1}pc = p^m c = 0$. Therefore, ka is not a generator of the cyclic group $\langle a \rangle$ since that group has p^m elements. Hence, $\text{GCD}(p^m, k) \neq 1$, and so p divides k . Let $k = pj$. Then $pb = ka = pji$. Let $b = c - ja$. Then $pb = 0$, but $b \notin \langle a \rangle$ as $c = b + ka \notin \langle a \rangle$. Q.E.D.

Proof. Let $|G| = p^n$ and $\text{ord } a = p^m$ with $m < n$.

We'll prove the lemma by induction. Assume it is valid for all groups of order less than p^n . Let b be an element in the complement of $\langle a \rangle$ of order p shown to exist in the previous lemma. Since $\text{ord } b = p$ and $b \notin \langle a \rangle$, therefore $\langle a \rangle \cap \langle b \rangle = 0$.

We'll reduce modulo $\langle b \rangle$ to a smaller p -primary group $G/\langle b \rangle$ where we can use the inductive hypothesis, then bring the results back up to G .

First, we'll show that $a + \langle b \rangle$, which is the image of a in $G/\langle b \rangle$, has the same order that a does in G , namely p^m , which implies that $a + \langle b \rangle$ is an element of maximal order in the group $G/\langle b \rangle$. Suppose $\text{ord}(a + \langle b \rangle) < p^m$. Then $p^{m-1}(a + \langle b \rangle)$ is the 0 element of $G/\langle b \rangle$, in other words, $p^{m-1}a \in \langle b \rangle$. But $p^{m-1}a \in \langle a \rangle$, and the intersection of $\langle a \rangle$ and $\langle b \rangle$ is trivial. Therefore, $p^{m-1}a = 0$ which contradicts $\text{ord } a = p^m$.

We now know $a + \langle b \rangle$ is an element of maximal order in the group $G/\langle b \rangle$, so we can apply the inductive hypothesis to conclude that $G/\langle b \rangle$ is the direct sum of the cyclic subgroup generated by $a + \langle b \rangle$ and another subgroup $K/\langle b \rangle$. Note that by the correspondence theorem, every subgroup of a quotient group $G/\langle b \rangle$ is the image of a group in G , so we may take K to be a subgroup of G .

We'll show that $G = \langle a \rangle \oplus K$ by showing that (1) $\langle a \rangle \cap K = 0$, and (2) $\langle a \rangle K = G$.

(1). If $x \in \langle a \rangle \cap K$, then its image $x + \langle b \rangle$ in the quotient group $G/\langle b \rangle$ lies in both the cyclic subgroup generated by $a + \langle b \rangle$ and $K/\langle b \rangle$. But their intersection is the 0 element in $G/\langle b \rangle$, therefore $x \in \langle b \rangle$. Since $x \in \langle a \rangle$ also, and $x \in \langle a \rangle \cap \langle b \rangle$ is trivial, therefore $x = 0$.

(2). We can show $\langle a \rangle K$ is all of G by a counting argument. We know that the order of $G/\langle b \rangle$ is the product of the order of the cyclic subgroup generated by $a + \langle b \rangle$ and the order of $K/\langle b \rangle$, the order of G is p times the order of $G/\langle b \rangle$, the order of $\langle a \rangle$ is the same as the order of the cyclic subgroup generated by $a + \langle b \rangle$, and the order of K is p times the order of

$K\langle b \rangle$. Therefore, the order of G equals the product of the order of $\langle a \rangle$ and the order of K . Thus $\langle a \rangle K = G$. Q.E.D.

You can prove the first statement of following theorem by induction using the lemma we just proved, then apply the primary decomposition theorem for the second statement. This is the existence half of the theorem we want. We'll still need some kind of uniqueness of the terms in the direct sum.

Theorem 4.68. A p -primary group is a direct sum of cyclic groups whose orders are powers of p . A finite Abelian group is the direct sum of cyclic groups.

There are a couple of ways to describe the uniqueness of the terms. Since we've been using cyclic groups whose orders are prime powers, let's stick to that.

There's a concept we'll need in the following lemma. If G is an Abelian group and p an integer, then the subset $G^p = \{x \mid px = 0\}$ is a subgroup of G . In fact, it's just the kernel of the group homomorphism $G \rightarrow G$ that maps x to px .

Exercise 88. Show that it is, indeed, a group homomorphism.

Lemma 4.69. Suppose that G is a p -primary group that can be written as a direct sum of nontrivial cyclic subgroups in two ways

$$G = H_1 \oplus H_2 \oplus \cdots \oplus H_m = K_1 \oplus K_2 \oplus \cdots \oplus K_n$$

where $|H_1| \geq |H_2| \geq \cdots \geq |H_m|$ and $|K_1| \geq |K_2| \geq \cdots \geq |K_n|$. Then $m = n$ and for each i , $|H_i| = |K_i|$.

Proof. Outline. By induction on the order of G . First verify that

$$G^p = H_1^p \oplus H_2^p \oplus \cdots \oplus H_m^p = K_1^p \oplus K_2^p \oplus \cdots \oplus K_n^p.$$

If any of the groups H_i^p or K_j^p are trivial, then drop them to get

$$G^p = H_1^p \oplus H_2^p \oplus \cdots \oplus H_{m'}^p = K_1^p \oplus K_2^p \oplus \cdots \oplus K_{n'}^p$$

to get two direct sums of nontrivial cyclic subgroups. By induction, $m' = n'$ and for each $i \leq m'$, $|H_i^p| = |K_i^p|$. Since $|H_i| = p|H_i^p|$ and $|K_i| = p|K_i^p|$, therefore $|H_i| = |K_i|$ for each $i \leq m'$. Finish with a counting argument to show that the number of trivial groups that were dropped is the same for the H 's as for the K 's. They're the subgroups H_i and K_i of order n . Q.E.D.

Putting the last theorem and lemma together, we have the following theorem.

Theorem 4.70 (Fundamental theorem of finite Abelian groups). A finite Abelian group is the direct sum of cyclic groups whose orders are prime powers. The number of terms in the direct sum and the orders of the cyclic groups are determined by the group.

Appendices

Appendix A

Background mathematics

A.1 Logic and proofs

Theorems. Logic and proofs are at the heart of mathematics. A statement will not be accepted by a mathematician if there's no proof of it. A theorem is a statement that has an accompanying proof.

If a statement is suspected to be true, but there's no proof yet, then it will be called a conjecture. Sometimes someone will supply the conjecture with a proof, then it becomes a theorem; sometimes a counterexample to the conjecture is discovered so it fails to be a theorem.

A typical theorem begins with the word “Theorem” followed by the statement of the theorem. That statement usually doesn't have much mathematical symbolism or variables, but it's written as much as possible in English sentences. After the proof is complete, it's ended with Q.E.D. (“Quod Erat Demonstrandum”, Latin for “that which was to be shown”) or some special symbol like a box \square .

Corollaries and lemmas are also theorems. A corollary is a theorem that follows quite easily from the preceding theorem. Sometimes the proof of a corollary is omitted and left to the reader to provide.

A lemma is a theorem that precedes another theorem. Lemmas are often technical and of little interest in themselves, but are necessary for the theorems which follows them. Sometimes a complicated proof will be split up and parts declared as lemmas. That makes it easier to understand the logical flow of the proof.

Some standard symbols seen in proofs. There are a whole lot of symbols and abbreviations that are used in proofs. Some are listed in table A.1. Although these are fairly standard, sometimes other symbols are used instead. They are used a lot when writing mathematics on a blackboard to save time. They're not as common in textbooks.

Besides these symbols, the symbol \because stands for “since”, and the symbol \therefore stands for “therefore”. They rarely appear in textbooks, but often on blackboards.

An example of universal quantification is the expression $\forall x, (x > 2 \Rightarrow x^2 > 4)$ which means for all x , if x is greater than 2, then x^2 is greater than 4. Typically conditions like $x > 2$ after universal quantifiers are included in the quantifier so that the implication doesn't have to be expressed separately. That last expression can be abbreviated as $\forall x > 2, x^2 > 4$.

Operation, symbol	Read As	Explanation
Conjunction, \wedge	and	The statement $A \wedge B$ is true if A and B are both true; else it is false.
Disjunction, \vee	(inclusive) or	The statement $A \vee B$ is true if A or B (or both) are true; if both are false, the statement is false.
Negation, \neg	not	The statement $\neg A$ is true just when A is false
Implication, \Rightarrow	implies; if... then	$A \Rightarrow B$ means if A is true, then B is also true; if A is false, then nothing is said about B .
Bi-implication, \Leftrightarrow	“iff”, if and only if	$A \Leftrightarrow B$ means <i>both</i> $A \Rightarrow B$ and $B \Rightarrow A$.
Universal quantification, \forall	for all; for any; for each	when it’s true universally
Existential quantification, \exists	there exists; there is an	when there’s at least one
Unique existential quantification, $\exists!$	there exists a unique	when there is exactly one

Table A.1: Standard logical symbols

An example of existential quantification is the expression $\exists x, (x > 1 \wedge x^2 = 4)$ which means there is an x such that x is greater than 1 and $x^2 = 4$. Typically conditions like $x > 1$ after existential quantifiers are included in the quantifier so that the conjunction doesn’t have to be expressed separately. That last expression can be abbreviated as $\exists x > 1, x^2 = 4$.

A.2 Sets

Just a little bit about sets. We’ll use the language of sets throughout the course, but we’re not using much of set theory. This note just collects the background that you need to know about sets in one place

A.2.1 Basic set theory

A set itself is just supposed to be something that has elements. It doesn’t have to have any structure but just have elements. The elements can be anything, but usually they’ll be things of the same kind.

If you’ve only got one set, however, there’s no need to even mention sets. It’s when several sets are under consideration that the language of sets becomes useful.

There are ways to construct new sets, too, and these constructions are important. The most important of these is a way to select some of the elements in a set to form another set, a subset of the first.

Examples. Let’s start with sets of numbers. There are ways of constructing these sets, but let’s not deal with that now. Let’s assume that we already have these sets.

The natural numbers. These are the counting numbers, that is, whole nonnegative numbers. That means we'll include 0 as a natural number. (Sometimes 0 isn't included.) There is a structure on \mathbf{N} , namely there are operations of addition, subtraction, etc., but as a set, it's just the numbers. You'll often see \mathbf{N} defined as

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

which is read as “ \mathbf{N} is the set whose elements are 0, 1, 2, 3, and so forth.” That's just an informal way of describing what \mathbf{N} is. A complete description couldn't get away with “and so forth.” If you want to see all of what “and so forth” entails, you can read Dedekind's 1888 paper *Was sind und was sollen die Zahlen?* and Joyce's comments on it. In that article Dedekind starts off developing set theory and ends up with the natural numbers.

The real numbers. These include all positive numbers, negative numbers, and 0. Besides the natural numbers, their negations and 0 are included, fractions like $\frac{22}{7}$, algebraic numbers like $\sqrt{5}$, and transcendental numbers like π and e . If a number can be named decimally with infinitely many digits, then it's a real number. We'll use \mathbf{R} to denote the set of all real numbers. Like \mathbf{N} , \mathbf{R} has lots of operations and functions associated with it, but treated as a set, all it has is its elements, the real numbers.

Note that \mathbf{N} is a subset of \mathbf{R} since every natural number is a real number.

Elements and membership. The standard notation to say an element x is a member of a set S is $x \in S$. The \in symbol varies a bit. Sometimes it appears as an epsilon ϵ or ε or \mathcal{E} . Read $x \in S$ as “ x is an element of S ,” or as “ x belongs to S , or more simply “ x is in S .”

It's negation is the symbol \notin . So, for example $\sqrt{5} \in \mathbf{R}$, but $\sqrt{5} \notin \mathbf{N}$.

As mentioned above, sets are completely determined by their elements, so two sets are equal if they have exactly the same elements.

$$S = T \text{ if and only if (1) for all } x \in S, x \in T, \text{ and (2) for all } x \in T, x \in S.$$

The two halves of the condition on the right lead to the concept of subset.

Subsets. If you have a set and a language to talk about elements in that set, then you can form subsets of that set by properties of elements in that language.

For instance, we have arithmetic on \mathbf{R} , so solutions to equations are subsets of \mathbf{R} . The solutions to the equation $x^3 = x$ are 0, 1, and -1 . We can describe its solution set using the notation

$$S = \{x \in \mathbf{R} \mid x^3 = x\}$$

which is read as “ S is the set of x in \mathbf{R} such that $x^3 = x$.” We could also describe that set by listing its elements, $S = \{0, 1, -1\}$. When you name a set by listing its elements, the order that you name them doesn't matter. We could have also written $S = \{-1, 0, 1\}$ for the same set. This set S is a subset of \mathbf{R} .

A set S is a subset of a set T if every element of S is also an element of T , that is

$$S \subseteq T \text{ if and only if for all } x \in S, x \in T.$$

Read $S \subseteq T$ as “ S is a subset of T .”

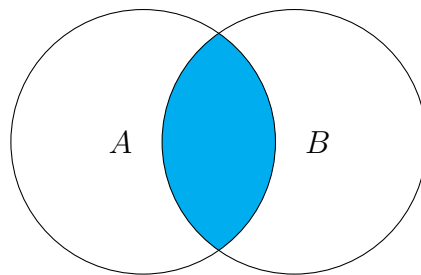
Note that $S = T$ if and only if $S \subseteq T$ and $T \subseteq S$.

There are a couple of notations for subsets. We'll use the notation $A \subseteq S$ to say that A is a subset of S . We allow $S \subseteq S$, that is, we consider a set S to be a subset of itself. If a subset A doesn't include all the elements of S , then A is called a *proper* subset of S . The only subset of S that's not a proper subset is S itself. We'll use the notation $A \subset S$ to indicate that A is a proper subset of S .

(Warning. There's an alternate notational convention for subsets. In that notation $A \subset S$ means A is any subset of S , while $A \subsetneq S$ means A is a proper subset of S . I prefer the notation we're using because it's analogous to the notations \leq for less than or equal, and $<$ for less than.)

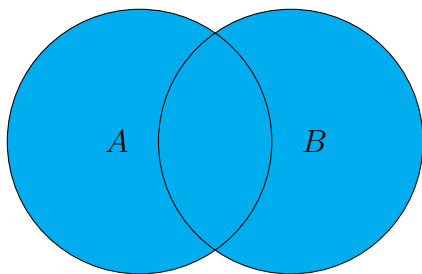
Operations on subsets. Frequently you deal with several subsets of a set, and there are operations of intersection, union, and difference that describe new subsets in terms of previously known subsets.

The intersection $A \cap B$ of two subsets A and B of a given set S is the subset of S that includes all the elements that are in both A and B , as shown in the Venn diagram below. (It's interesting that Venn called them Euler circles as Euler had used them earlier, but Leibniz had also used them, and Ramon Llull (Raymond Lully) in the 13th century.) Read $A \cap B$ as "the intersection of A and B " or as " A intersect B ." Note that the operation of intersection is associative and commutative.



$$A \cap B = \{x \in S \mid x \in A \text{ and } x \in B\}.$$

Two sets A and B are said to be *disjoint* if their union is empty, $A \cap B = \emptyset$. Several sets are said to be *pairwise* disjoint if each pair of those sets are disjoint.



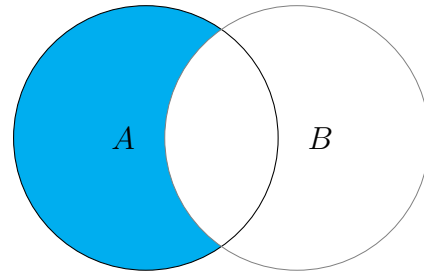
$$A \cup B = \{x \in S \mid x \in A \text{ or } x \in B\}.$$

The union $A \cup B$ of two subsets A and B of a given set S is the subset of S that includes all the elements that are in A or in B or in both. Read $A \cup B$ as "the union of A and B " or as " A union B ." Like intersection, the operation of union is also associative and commutative. It is usual in mathematics to take the word "or" to mean an inclusive or. It implicitly includes "or both."

Intersection and union each distribute over the other:

$$\begin{aligned} (A \cap B) \cup C &= (A \cup C) \cap (B \cup C) \\ (A \cup B) \cap C &= (A \cap C) \cup (B \cap C) \end{aligned}$$

The difference $A - B$ of two subsets A and B of a given set S is the subset of S that includes all the elements that are in A but not in B .



$$A - B = \{x \in S \mid x \in A \text{ and } x \notin B\}$$

There's also the complement of a subset A of a set S . The complement is just $S - A$, all the elements of S that aren't in A . When the set S is understood, the complement of A often is denoted more simply as either A^c , \bar{A} , or A' rather than $S - A$. I prefer the notation A^c .

These operations satisfy lots of identities. I'll just name a couple of important ones.

De Morgan's laws describe a duality between intersection and union. They can be written as

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

Unions and intersections sometimes are taken of many subsets, even infinitely many. Suppose that A_1, A_2, \dots, A_n are subsets of S . The intersection of all of them can be written in an indexed notation as

$$\bigcap_{i=1}^n A_i = A_1 \cap A_2 \cap \dots \cap A_n$$

and their union as

$$\bigcup_{i=1}^n A_i = A_1 \cup A_2 \cup \dots \cup A_n.$$

And when there are infinitely many, $A_1, A_2, \dots, A_n, \dots$, as

$$\bigcap_{i=1}^{\infty} A_i = \{x \in S \mid x \in A_i \text{ for all } i\}$$

and their union as

$$\bigcup_{i=1}^{\infty} A_i = \{x \in S \mid x \in A_i \text{ for at least one } i\}.$$

DeMorgan's laws and the distributivity laws also apply to indexed intersections and unions.

$$\begin{aligned} \left(\bigcap_{i=1}^n A_i \right)^c &= \bigcup_{i=1}^n A_i^c \\ \left(\bigcup_{i=1}^n A_i \right)^c &= \bigcap_{i=1}^n A_i^c \\ \left(\bigcap_{i=1}^n A_i \right) \cup C &= \bigcap_{i=1}^n (A_i \cup C) \\ \left(\bigcup_{i=1}^n A_i \right) \cap C &= \bigcup_{i=1}^n (A_i \cap C) \end{aligned}$$

Partitions. A set S is said to be *partitioned* into subsets A_1, A_2, \dots, A_n when each element of S belongs to exactly one of the subsets A_1, A_2, \dots, A_n . That's logically equivalent to saying that S is the disjoint union of the A_1, A_2, \dots, A_n .

When you have a partition A_1, A_2, \dots, A_n of a set S like that, it induces a partition $E \cap A_1, E \cap A_2, \dots, E \cap A_n$ on each subset E of S . Each element of E belongs to exactly one of its subsets $E \cap A_1, E \cap A_2, \dots, E \cap A_n$.

Products of sets. So far we've looked at creating sets within set. There are some operations on sets that create bigger sets, the most important being creating products of sets. These depend on the concept of ordered pairs of elements. The notation for ordered pair (a, b) of two elements extends the usual notation we use for coordinates in the xy -plane. The important property of ordered pairs is that two ordered pairs are equal if and only if they have the same first and second coordinates:

$$(a, b) = (c, d) \text{ iff } a = c \text{ and } b = d.$$

The product of two sets S and T consists of all the ordered pairs where the first element comes from S and the second element comes from T :

$$S \times T = \{(a, b) \mid a \in S \text{ and } b \in T\}.$$

Thus, the usual xy -plane is $\mathbf{R} \times \mathbf{R}$, usually denoted \mathbf{R}^2 .

Besides binary products $S \times T$, you can analogously define ternary products $S \times T \times U$ in terms of triples (a, b, c) where $a \in S$, $b \in T$, and $c \in U$, and higher products, too.

Sets of subsets; power sets. Another way to create bigger sets is to form sets of subsets. If you collect all the subsets of a given set S into a set, then the set of all those subsets is called the *power set* of S , denoted $\wp(S)$ or sometimes 2^S .

For example, let S be a set with 3 elements, $S = \{a, b, c\}$. Then S has eight subsets. There are three singleton subsets, that is, subsets having exactly one element, namely $\{a\}$, $\{b\}$, and $\{c\}$. There are three subsets having exactly two elements, namely $\{a, b\}$, $\{a, c\}$, and $\{b, c\}$. There's one subset having all three elements, namely S itself. And there's one subset that has no elements. You could denote it $\{\}$, but it's always denoted \emptyset and called the *empty set* or *null set*. Thus, the power set of S has eight elements

$$\wp(S) = \{\emptyset, \{a\}, \{b\}, \{c\}, \{a, b\}, \{a, c\}, \{b, c\}, S\}.$$

Cardinality, countable versus uncountable sets. The *cardinality* of a set S is the number of elements in it, denoted $|S|$. So, for example, if $S = \{a, b, c\}$, then $|S| = 3$, and $|\wp(S)| = 2^3 = 8$.

Some sets are infinite, so their cardinality is not a finite number. A more careful definition is needed. Two sets S and T are said to have the *same cardinality* if there is a one-to-one correspondence of their elements. That means that there is some function $f : S \rightarrow T$ which is injective (also called one-to-one) and surjective (also called onto). A function which is both injective and surjective is called a *bijection*. For a bijection $f : S \rightarrow T$, the inverse function $f^{-1} : T \rightarrow S$ is also a bijection. The notation $|S| = |T|$ indicates S and T have the same cardinality.

If there is an injection $S \rightarrow T$, then the cardinality of S is less than or equal to that of T , written $|S| \leq |T|$. It is evident that \leq is a transitive relation on cardinalities. The Schröder-Bernstein theorem states that if there are injections both ways between S and T , then they have the same cardinality. Thus, \leq is a partial order on cardinalities.

The notation $|S| < |T|$ means $|S| \leq |T|$ but not $|S| = |T|$.

As Georg Cantor (1845–1918) discovered, not all infinite sets have the same cardinality. Some infinite sets are bigger than others. Using his famous diagonal proof, he proved that for any set, even if it's infinite, $|S| < |\mathcal{P}(S)|$.

The smallest size an infinite set can be is that of the natural numbers \mathbf{N} . A set that has the same cardinality as \mathbf{N} is called a *countably infinite* set. An infinite set that doesn't have the same cardinality as \mathbf{N} is called an *uncountable* set. The set of real numbers \mathbf{R} is uncountable.

Finite sets are also said to be countable. Thus, a set is *countable* if it's either finite or countably infinite.

A.2.2 Functions and relations

A function f is associated to a pair of sets, a *domain* S and a *codomain* T . The usual notations for that are $f : S \rightarrow T$ and $S \xrightarrow{f} T$. In order to be a function, each element $x \in S$ must be associated to a particular element of T , denoted $f(x)$.

The *graph* of a function f is a subset of the product $S \times T$, namely, the set $F = \{(x, y) \in S \times T \mid y = f(x)\}$.

Two functions are said to be the same if they have the same graph, so the graph characterizes the function. Frequently, textbooks define a function $f : S \rightarrow T$ as its graph, that is, a subset F of $S \times T$ such that for all $x \in X$, there is a unique $y \in T$ such that $(x, y) \in F$.

When $f : S \rightarrow T$, it is said that f *maps* S to T , and that f *maps* x to $f(x)$. This element $f(x)$ is called the *image* of x under f . The mapping of x to $f(x)$ is denoted $x \mapsto f(x)$.

The concept of image is extended to subsets of the domain. If $A \subseteq S$, then f maps A to the set $f(A) = \{f(x) \mid x \in A\}$, called the *image* of A under f .

Another related concept is that of preimage, also called inverse image. If B is a subset of the codomain T , then the preimage of B under f is the set $f^{-1}(B) = \{x \in A \mid f(x) \in B\}$.

Composition. If $f : S \rightarrow T$ and $g : T \rightarrow U$, then the composition $g \circ f : S \rightarrow U$ is defined by $(g \circ f)(x) = g(f(x))$.

Composition is associative. $(h \circ g) \circ f = h \circ (g \circ f)$. Since composition is associative, parentheses are not necessary when composing three or more functions.

For each set S there is an identity function $1_S : S \rightarrow S$ which maps every element in S to itself, $1_S(x) = x$. The identity functions act as units for composition. If $f : S \rightarrow T$, then $1_T \circ f = f$ and $f = f \circ 1_S$.

Injections, surjections, and bijections. These are words that describe certain functions $f : S \rightarrow T$ from one set to another. An *injection*, also called a *one-to-one function* is a function that maps distinct elements to distinct elements, that is, if $x \neq y$, then $f(x) \neq f(y)$. Equivalently, if $f(x) = f(y)$ then $x = y$. If S is a subset of T , then there is a natural injection $\iota : S \rightarrow T$, called the *inclusion function*, defined by $\iota(x) = x$.

A *surjection*, also called an *onto function* is one that includes all of T in its image, that is, if $y \in T$, then there is an $x \in S$ such that $f(x) = y$.

A *bijection*, also called a *one-to-one correspondence*, is a function that is simultaneously injective and bijective. Another way to describe a bijection is to say that there is an inverse function $g : T \rightarrow S$ so that the composition $g \circ f : S \rightarrow S$ is the identity function on S while $f \circ g : T \rightarrow T$ is the identity function on T . The usual notation for the function inverse to f is f^{-1} . In this situation f and g are inverse to each other, that is, if g is f^{-1} , then f is g^{-1} . Thus, $(f^{-1})^{-1} = f$.

Relations. Relations include functions, but are more general. A binary *relation* $R : S \rightarrow T$ doesn't have to associate each element of S to exactly one element of T . It can associate an element of S to any number of elements in T including the possibility of no elements in T at all. In other words, a relation $R : S \rightarrow T$ is determined by an arbitrary subset of $S \times T$.

The most useful relations are those that have special properties. The next section discusses equivalence relations. A typical equivalence relation is congruence modulo n . Order relations are discussed in section A.3. A typical order relation is \leq on numbers.

A.2.3 Equivalence relations

There are various symbols used for equivalence relations, such as \cong , \equiv , \approx , \asymp , \simeq , \sim , and so forth. We'll use \equiv for a generic equivalence relation.

Definition A.1 (Equivalence relation). An *equivalence relation* \equiv on a set S is a relation that is reflexive, symmetric, and transitive.

A relation on a set S may be identified with a subset of the product $S \times S$. For an equivalence relation \equiv , this means $x \equiv y$ corresponds to the statement that the ordered pair (x, y) is an element of that subset.

Reflexivity: For all x , $x \equiv x$.

Symmetry: For all x and y , $x \equiv y$ implies $y \equiv x$.

Transitivity: For all x , y , and z , $x \equiv y$ and $y \equiv z$ implies $x \equiv z$.

Equivalence classes and partitions of sets. An equivalence relation on a set determines a partition on that set, and conversely, as we'll see presently.

Definition A.2 (Equivalence class). Given an equivalence relation on a set, an *equivalence class* of an element x , denoted $[x]$, is the set of all elements equivalent to x ,

$$[x] = \{y \mid y \equiv x\}.$$

You can easily show the several properties of equivalence classes.

Theorem A.3. If \equiv is an equivalence relation on a set S , then the following four statements are equivalent

1. $x \equiv y$.
2. $[x] = [y]$.

3. $x \in [y]$.
4. $[x] \cap [y] \neq \emptyset$.

Furthermore, for each $x \in S$, there is exactly one equivalence class containing x , namely, $[x]$.

Definition A.4 (Partition of a set). A *partition* of a set S is a collection of nonempty subsets, called *parts*, of S which are pairwise disjoint and whose union is all of S . Thus, each element of S belongs to exactly one of the parts.

The above theorem shows that the equivalence classes form a partition. The converse is also true as you can easily show.

Theorem A.5. For each equivalence class on a set, the equivalence classes partition the set. Conversely, a partition of a set determines an equivalence relation where two elements are equivalent if they're in the same part.

The set of equivalence classes is sometimes denoted S/\equiv , and it's sometimes called a quotient set. Using equivalence classes to construct new sets of things is a common practice in mathematics and especially in algebra.

Keep in mind that you can always name an element of S/\equiv by naming an element of S , but two elements x and y of S will name the same element of S/\equiv , that is, $[x] = [y]$, if $x \equiv y$.

The function $\gamma : S \rightarrow S/\equiv$ defined by $\gamma(x) = [x]$ is called a *projection*, or the *canonical function*, from the set to its quotient set.

A.2.4 Axioms of set theory

Although the axioms of set theory don't play an important role in an introductory course in modern algebra, occasionally they may be useful. Here is a summary of axioms of Zermelo-Fraenkel set theory, abbreviated ZF set theory.

Axiom of extensionality . This is the axiom that says two sets are the same if they have the same elements.

$$\forall A, \forall B, (\forall x, (x \in A \Leftrightarrow x \in B) \iff A = B).$$

Axiom of separation . This axiom is also called the axiom of specification. It says if you have a predicate φ on sets and a given set A , then there is a subset B of A on which that predicate holds.

$$\forall A, \exists B, \forall x, (x \in B \iff x \in A \wedge \varphi(x)).$$

It's an axiom schema rather than a single axiom because a different axiom is needed for each predicate φ .

This axiom allows the creation of smaller sets from a given set. For example, if $A = \mathbf{R}$ is the set of real numbers, by the axiom of separation there is a set B such that the elements of B are the real numbers that satisfy the equation $x^3 - 3x = 1$. Here, the predicate φ at x , written above as $\varphi(x)$ is that equation. The axiom of separation is the justification for the "set building" notation $B = \{x \in \mathbf{R} \mid x^3 - 3x = 1\}$.

Axiom of pairing. The axiom of pairing allows the creation of a set containing two elements (or one if they're the same element).

$$\forall x, \forall y, \exists A, (z \in A \iff z = x \vee z = y).$$

The set A is usually denoted $\{x, y\}$.

If it happens that $x = y$, then A only has one element instead of two since $\{x, x\} = \{x\}$. An set with only one element is called a *singleton set*, or just a *singleton*.

Axiom of union. Given a set A of sets, this says the union C of the sets in A is a set.

$$\forall A, \exists C, \forall x, (x \in C \iff \exists B, (x \in B \wedge B \in A)).$$

The usual notation for C is $\bigcup A$, or $\bigcup_{B \in A} B$.

When A is the pair $\{D, E\}$ then $\bigcup_{B \in A} B$ is the pairwise union $D \cup E$.

There doesn't need to be an axiom for intersections or for relative compliments because intersections and relative complements can be proved from the axiom of separation.

Axiom of powersets. It says given a set A , there is a set which contains all the subsets of A .

$$\forall A, \exists B, \forall C, (C \in B \iff C \subseteq A).$$

One common notation for the powerset B of A is $\mathcal{P}(A)$.

Axiom of infinity. So far, there are no axioms that say there are any sets at all. This axiom says that there is an infinite set which contains the emptyset \emptyset , so among other things, it says the emptyset exists. When studying the theory of finite sets, the axiom of infinity is not included, but an explicit axiom is needed to say the emptyset exists.

Define $S(A)$ to denote $A \cup \{A\}$ where A is a set. $S(A)$ is called the *successor* of A . The axiom of pairing says that if A is a set, then so is $\{A\}$, and the axiom of union then implies that $A \cup \{A\}$ is a set.

The axiom of infinity says that there is a set B that has \emptyset as an element and is closed under S .

$$\exists B, (\emptyset \in B \wedge \forall y \in B, S(y) \in B).$$

Along with the axiom of regularity, this axiom implies that there is at least one infinite set. With the other axioms, it can be shown that there is a smallest such set. That smallest set is a model for the set of natural numbers \mathbf{N} . In that model, the emptyset \emptyset acts as 0, it's successor $S(\emptyset)$ acts as 1, $S(S(\emptyset))$ acts as 2, and so forth.

Axiom of regularity. This axioms is also called the axiom of foundation. It is a technical axiom that says that given a nonempty set A , there is an element of A which is disjoint from it.

$$\forall A \neq \emptyset, \exists x \in A, \forall y \in x, y \notin A.$$

The axiom of regularity implies that no set is an element of itself, nor is there a finite cycle of memberships where $A_1 \in A_2 \in \dots \in A_n \in A_1$. Furthermore, there is no infinite descending memberships $\dots \in A_n \in \dots \in A_2 \in A_1$. One of the main reasons for the axiom of regularity is to develop the theory of ordinals.

Axiom of replacement. Like the axiom of separation, this is another axiom schema. This technical axiom creates images of functions described by predicates. A predicate φ describes a function if for all x , there exists a unique y such that $\varphi(x, y)$. In that case, a function symbol like f is used so that $f(x) = y$ expresses $\varphi(x, y)$. (For this axiom, the predicate can have other arguments that won't be mentioned explicitly.)

$$\forall A, ((\forall x \in A, \exists! y, \varphi(x, y)) \Rightarrow \exists B, \forall x \in A, \exists y \in B, \varphi(x, y)).$$

The B in the axiom is usually denoted $f(A)$, the image of A under f .

Axiom of choice. The axiom of choice is not part of ZF set theory, but when it's included, the set theory is denoted ZFC set theory, Zermelo-Fraenkel set theory with the axiom of choice. This axiom is discussed in more detail in the section [A.4](#).

Von Neumann–Bernays–Gödel set theory (NBG). This is an extension of ZF that includes proper classes. Whereas sets can be elements of other sets and proper classes, proper classes cannot be elements. NBG is a conservative extension of ZF in that sense that any theorem not mentioning classes and provable in one theory can be proved in the other. NBG makes it possible to talk about things like the class of all sets, or the class of all groups, etc.

A.3 Ordered structures

Several mathematical structures are defined in terms of an order relation. These order relations have something in common with the order relation \leq “less than or equal” on the real numbers. Many of them are not total orders like \leq , but only partial orders. Having fewer nice properties than \leq , however, can make them more interesting.

In particular, we'll look at partial orders, lattices, and Boolean algebras.

A.3.1 Partial orders and posets.

You're familiar with the order \leq on real numbers. It's a binary relation with the following four properties.

1. Reflexivity: for all x , $x \leq x$.
2. Anti-symmetry: for all x and y , if $x \leq y$ and $y \leq x$, then $x = y$.
3. Transitivity: for all x , y , and z , if $x \leq y$ and $y \leq z$, then $x \leq z$.
4. Totality: for all x and y , either $x \leq y$ or $y \leq x$ (or both in which case $x = y$).

There are other useful binary relations in mathematics with either those four properties or at least the first three. Although sometimes such binary relations are denoted with the same \leq sign, frequently a similar but visually distinct sign such as \preceq is used. Both are read “less than or equal to”. Of course, there's also a greater than or equal to, written \succeq and defined by $x \succeq y$ if and only if $y \preceq x$.

Definition A.6 (Total order). A *total order*, also called a *linear order*, on a set is a binary relation having the four properties: reflexivity, anti-symmetry, transitivity, and totality.

A set with a specified total order is called a *totally ordered set* or a *chain*.

The strict form \prec of a total order (or a partial order defined below) \preceq is defined by

$$x \prec y \text{ if and only if } x \preceq y \text{ and } x \neq y.$$

A useful weakening of total orders is what is called a partial order.

Definition A.7 (Partial order). A *partial order* on a set is a binary relation having the first three of those properties: reflexivity, anti-symmetry, and transitivity.

A set with a specified partial order is called a *partially ordered set* or *poset* for short.

Two elements x and y in a partially ordered set are said to be *comparable* if either $x \preceq y$ or $y \preceq x$. Otherwise they're *incomparable*.

Example A.8. The positive integers are partially ordered by divisibility. Divisibility is reflexive since $x|x$; it's anti-symmetric since if $x|y$ and $y|x$, then $x = y$, and it's transitive since if $x|y$ and $y|z$, then $x|z$.

Divisibility is not a partial order on all integers since $2|-2$ and $-2|2$ but $2 \neq -2$. It is, however, a pre-order. A *pre-order* is reflexive and transitive but need not be anti-symmetric.

Example A.9. Any collection \mathcal{T} of subsets of a set S is a partially ordered set where the binary relation is \subseteq . In particular, the power set $\mathcal{P}(S)$ consisting of all the subsets of S is partially ordered.

Hasse diagrams. A partially ordered set can be described with a kind of a graph called a Hasse diagram. The elements of the set are the vertices of the graph, and the edges indicate which elements are less than or equal to which other elements. If $a \prec b$, then an edge is drawn from a to b with the larger element above the smaller one. Transitivity of the order relation is assumed so that if $a \prec b \prec c$, then an edge doesn't have to be drawn between a and c .

Definition A.10. An *upper bound* of a subset S in a poset is any element in the poset which is greater than or equal to all elements in S . That element needn't be an element of the subset S . Likewise, a *lower bound* of S is an element that is less than or equal to all the elements in S .

A *least upper bound*, also called *supremum* of S is an upper bound of S which is less than or equal to all other upper bounds of S . It is denoted $\text{lub } S$ or $\text{sup } S$. Likewise, a *greatest lower bound*, also called *infimum* of S is an lower bound of S which is greater than or equal to all other lower bounds of S . It is denoted $\text{glb } S$ or $\text{inf } S$.

Least upper bounds and greatest lower bounds of subsets need not always exist.

Example A.11. With the usual ordering on the real numbers \mathbf{R} , both the open interval $(2, 3)$ and the closed interval $[2, 3]$ have the same least upper bound 3 and the same greatest lower bound 2.

With the usual ordering on the rational numbers \mathbf{Q} , the subset $S = \{x | x^2 = 2\}$ has neither a least upper bound nor a greatest lower bound since $\sqrt{2}$ and $-\sqrt{2}$ are not rational numbers.

Definition A.12 (Maximal and minimal elements). A *maximal element* in a partially ordered set is an element which is not less than or equal to any other element. A *minimal element* in a partially ordered set is an element which is not greater than or equal to any other element.

Maximal and minimal elements don't have to be unique. A partially ordered set can have more than one of each or none at all.

Definition A.13 (Meet and join). The *meet* of two elements a and b is the greatest lower bound of the set $\{a, b\}$. That is, it is an element x less than or equal to both a and b and greater than or equal to all other elements greater than or equal to both a and b . If that meet exists, it is denoted $a \wedge b$.

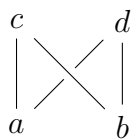
The *join* of two elements a and b in a partially ordered set is the least upper bound of the set $\{a, b\}$. That is, it is an element x greater than or equal to both a and b and less than or equal to all other elements greater than or equal to both a and b . If that join exists, it is denoted $a \vee b$.

Meets and joins aren't particularly interesting in totally ordered sets. In a totally ordered set, the meet of two elements is the minimum of the two while the join of two elements is the maximum of the two.

Example A.14. Consider the positive integers partially ordered by divisibility. The meet of two integers m and n is that number d which divides them both for which any other divisor of both divides d . In other words, a meet in this partially ordered set is the greatest common divisor.

Likewise, a join is the least common multiple.

Example A.15. Sometimes meets and joins don't exist in a partially ordered set. Consider the poset with four elements, $\{a, b, c, d\}$ where both a and b are less than both c and d .



The join $c \vee d$ doesn't exist since there is no upper bound for c and d . The join $a \vee b$ doesn't exist because there are two upper bounds for a and b , but no least upper bound. Likewise, the two meets $c \wedge d$ and $a \wedge b$ don't exist.

A.3.2 Lattices

Lattices are partially ordered sets that have meets and joins

Definition A.16 (Lattice). A *lattice* is a partially ordered set in which all meets and joins of two elements exist, has a smallest element (often denoted 0 or \perp) and a largest element (often denoted 1 or \top), in which the following identities hold.

Idempotency: $x = x \wedge x = x \vee x$.

Commutativity: $x \wedge y = y \wedge x$ and $x \vee y = y \vee x$.

Associativity: $(x \wedge y) \wedge z = x \wedge (y \wedge z)$ and $(x \vee y) \vee z = x \vee (y \vee z)$.

Absorption: $a \wedge (a \vee b) = a$ and $a \vee (a \wedge b) = a$.

Identity: $a \wedge 1 = a$ and $a \vee 0 = a$.

Lattices can be defined without reference to a partial order as the relation $a \leq b$ can be characterized in terms of meets and joins as in the following theorem.

Theorem A.17. The following three conditions are equivalent: $a \preceq b$, $a \wedge b = a$ and $a \vee b = b$.

Proof. First, suppose $a \preceq b$, then by definition, the meet of a and b is a while the join of a and b is b . Thus, the first condition in the statement of the theorem implies the other two.

Now suppose $a \wedge b = a$, since $a \wedge b \preceq b$, therefore $a \preceq b$. Thus the second condition implies the first. Similarly, the third condition implies the first. Q.E.D.

Since \preceq can be characterized in terms of \wedge and \vee , there is an alternate definition of lattice.

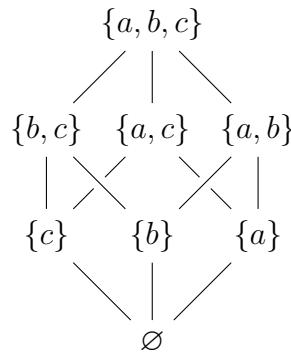
Definition A.18 (Alternate definition of lattice). A *lattice* is a set equipped with two binary operations \wedge and \vee and two constants 0 and 1 which satisfy the identities in the previous definition.

It follows from that definition that $a = a \wedge b$ if and only if $b = b \vee a$. (Proof: $a = a \wedge b$ implies $b = b \vee (b \wedge a) = (a \wedge b) \vee b = a \vee b$ by commutativity and absorption.)

The partial order can then be recovered by defining $a \preceq b$ if and only if $a \wedge b = a$ and $a \vee b = b$.

There are a couple other identities that follow from the definition, namely, $0 \wedge a = 0$ and $1 \vee a = 1$.

Example A.19. The powerset $\mathcal{P}(S)$ of a set S is a lattice. (It's actually a Boolean ring, discussed later.) Here's the Hasse diagram for $\mathcal{P}(\{a, b, c\})$.



The powerset of a set with four elements has 16 elements. It's a little harder to draw as a Hasse diagram which is displayed in figure A.1. The names of the subsets are abbreviated so that, for example, the subset $\{a, b, c\}$ is displayed as abc .

Modular and distributive lattices. Note that distributivity is not listed among the identities above. That's because it doesn't hold in all lattices. Another identity that doesn't hold in all lattices is modularity.

Definition A.20. A lattice is said to be *modular* if for all elements a , b , and c for which $a \preceq c$, it is the case that $a \vee (b \wedge c) = (a \vee b) \wedge c$.

A lattice is said to be *distributive* if for all elements a , b , and c , it is the case that $a \wedge (b \vee c) = (a \wedge b) \vee (a \wedge c)$ and $a \vee (b \wedge c) = (a \vee b) \wedge (a \vee c)$.

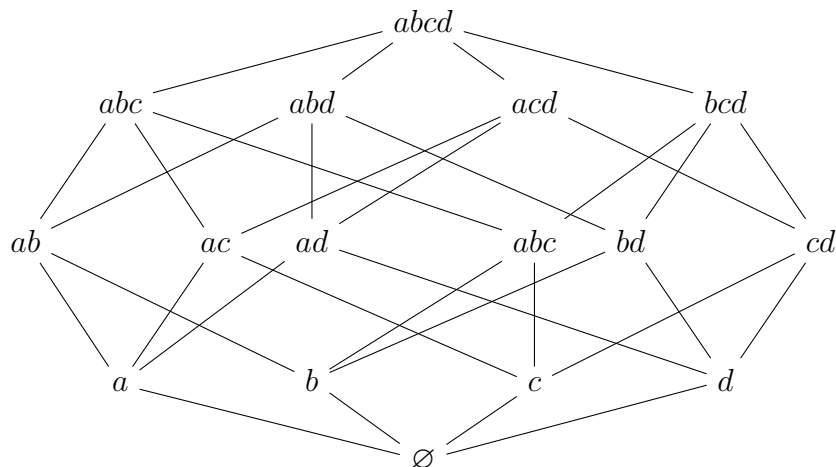


Figure A.1: Lattice of the Powerset of 4 elements

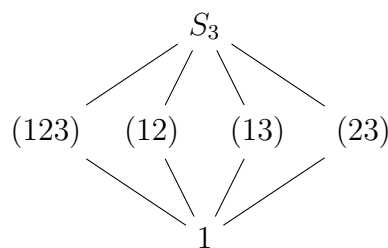
It can be shown that every distributive lattice is also modular, but there are modular lattices that are not distributive. It can also be shown that either one of the distributive identities imply the other.

The powerset $\mathcal{P}(S)$ lattice is a distributive lattice.

The subgroups of a group with inclusion as a partial order always form a modular lattice, but not always a distributive lattice.

Example A.21.

The symmetric group S_3 has four subgroups besides itself and the trivial subgroup. The subgroup generated by the permutation (123) has order 3 while the three subgroups generated by the three transpositions (12), (13), and (23) each have order 2. The lattice of subgroups is modular, but it's not distributive.



A.3.3 Boolean algebras.

A Boolean algebra is a distributive lattice with one more operation.

Definition A.22 (Boolean algebra). A *Boolean algebra* is a distributive lattice with a unary operation, called *complementation* or *negation*, denoted \neg satisfying the identities $a \vee \neg a = 1$ and $a \wedge \neg a = 0$.

Actually, not all the identities from boolean lattices are necessary for the definition since absorption can be shown from the rest. Other identities that follow from the definition include $\neg 0 = 1$, $\neg 1 = 0$, and $\neg \neg a = a$.

As described in section 3.3, Boolean algebras are the same thing as Boolean rings. The only difference is notational.

Truth values. The two-element Boolean algebra that consists only of 0 and 1 is used in logic. 0, or \perp is the truth value “false” while 1, or \top is the truth value “true”.

A.4 Axiom of choice

Given a collection of nonempty sets, the axiom of choice says that there is a function that chooses one element from each set.

This is an axiom of set theory. There are many axioms of set theory, most of which are fairly obvious and uncontroversial.

More precisely, the axiom of choice says that given any set S , there exists a “choice function” $\gamma : \mathcal{P}(S) - \emptyset \rightarrow S$ which chooses from any nonempty set $T \subseteq S$ an element $\gamma(T) \in T$.

In some sense, any theorem that relies on the axiom of choice is flawed since the axiom of choice is not constructive. So, for instance, after proving an ideal is a subideal of a maximal ideal, we won’t have any way to identify that maximal ideal.

Here’s a simple theorem that relies on the axiom of choice.

Theorem A.23. Let $f : A \rightarrow B$ be a surjective function between sets A and B . Then there exists $g : B \rightarrow A$ such that $f \circ g$ is the identity function on B .

Proof. Let γ be a choice function for A . Then g is the function

$$g(y) = \gamma(f^{-1}(y)) = \gamma(\{x \mid f(x) = y\}).$$

Since f is surjective, $f^{-1}(y)$ is not the empty set, so the choice function γ will choose some element x out of $f^{-1}(y)$ with $f(x) = y$. Q.E.D.

That theorem is actually logically equivalent to the axiom of choice, that is, the axiom of choice follows from it.

Independence of the axiom of choice. The axiom of choice is independent of the rest of the axioms of set theory. Gödel proved in 1938 that set theory with the axioms of choice added is as consistent as set theory, while Cohen in 1963 proved that set theory with the negation of the axiom of choice added is as consistent as set theory. In other words, the axiom of choice is independent of the rest of the axioms.

A.4.1 Zorn’s lemma

Although the axiom of choice is easy to state, it’s not usually easy to use. Zorn’s lemma, which is logically equivalent to the axiom of choice is hard to state, but easy to use. Another is the well-ordering principle.

This lemma is applied to a nonempty collection \mathcal{M} of subsets of a set S .

Section A.3.1 on partially-ordered sets defined a chain, upper bound, and maximal element. A chain in \mathcal{M} is a collection \mathcal{C} of elements of \mathcal{M} linearly ordered by subset inclusion. In other words, if A and B are elements of \mathcal{C} , either $A \subseteq B$ or $B \subseteq A$. An upper bound of \mathcal{C} is a subset B of S which contains all elements of \mathcal{C} . A maximal element B of \mathcal{M} is one not contained in any larger element of \mathcal{M} .

Zorn’s lemma. If every chain in \mathcal{M} has an upper bound in \mathcal{M} , then \mathcal{M} has a maximal element.

We won’t prove that the Axiom of Choice is equivalent to Zorn’s lemma because it would take too long.

A.4.2 Well-ordering principle

The most common form of the axioms of choice used in algebra is Zorn's lemma. Another that's sometimes used is the well-ordering principle.

Definition A.24 (Well-ordering). A partially ordered set is well ordered if every nonempty subset of it has a least element.

It follows from the definition that every well-ordering is totally ordered. Given two elements x and y , the subset $\{x, y\}$ has a smallest element, either x in which case $x \preceq y$, or y in which case $y \preceq x$.

Example A.25. The natural numbers \mathbf{N} is well-ordered by its usual ordering. The integers \mathbf{Z} is not well ordered by its usual ordering because the entire set doesn't have a smallest element. For the same reason \mathbf{R} is not well-ordered. The non-negative real numbers aren't well-ordered by its usual ordering because any open interval (a, b) doesn't have a least element.

Any subset of a well-ordered set is well-ordered by the same ordering.

Lexicographic ordering. The product $\mathbf{N} \times \mathbf{N}$ has a well-ordering called the lexicographic ordering. The ordered pair (a, b) is defined to be less than or equal to the ordered pair (c, d) if either $a = c$ and $b \leq d$ or $a \leq c$. Thus, the elements of $\mathbf{N} \times \mathbf{N}$ listed in increasing order are

$$(0, 0), (0, 1), (0, 2), \dots; (1, 0), (1, 1), (1, 2), \dots; (2, 0), (2, 1), (2, 2), \dots; \dots$$

More generally, if A and B are both well ordered, then the lexicographic order on $A \times B$ is a well-ordering.

Furthermore, finite products $A_0 \times A_2 \cdots A_2$ of well-ordered sets are well ordered by a lexicographic ordering.

The well-ordering principle . This principle states that every set has a well-ordering, that is, for each set, there is some well-ordering of that set.

The axiom of choice, Zorn's lemma, and the well-ordering principle can each be proved from the each other. Here's a proof that the well-ordering principle follows from Zorn's lemma.

Theorem A.26. The well-ordering principle follows from Zorn's lemma.

Proof. Let S be a set. Let \mathcal{W} be the set of well-orderings of subsets of S . Partially order \mathcal{W} so given subsets A and B both with will-orderings, define $A \preceq B$ if $A \subseteq B$ and the two orderings agree on A . In other words, the ordering on A extends to that on B .

To use Zorn's lemma, we need to show that every chain \mathcal{C} in \mathcal{W} has an upper bound. A chain \mathcal{C} consists of subsets A and B where if $A \preceq B$, the ordering of A is extended to B . The union of all these subsets is a set C which, when given the extended ordering, so lies in \mathcal{W} , is itself a well-ordered set that contains every subset $A \in \mathcal{C}$. Thus, every chain in \mathcal{W} has an upper bound.

By Zorn's lemma, \mathcal{W} has a maximal element M . This M is a well-ordered subset of S which cannot be extended (since it's maximal). If there were an element of $S - M$, the

ordering on M could be extended to a well-ordering by making that element less than every element in M . Therefore, there are no elements in $S - M$. Thus, $S = M$, and so S has a well-ordering. Q.E.D.

This principle implies that there is some well-ordering of the real numbers \mathbf{R} . It's not the usual order, of course, since the usual order does not well order \mathbf{R} . In fact, no particular well-ordering of \mathbf{R} can ever be described.

Index

- A_4 , 108, 112
- A_n , *see* Alternating group
- D_3 , 8
- D_5 , 108, 111, 119
- D_n , *see* Dihedral group
- $F_{21} = C_7 \rtimes C_3$, 133
- $GF(p^n)$, *see* Galois field
- $GL_2(\mathbf{R})$, 7
- $GL_n(R)$, *see* General linear group
- $PGL_n(F)$, *see* Projective linear group
- S^1 , *see* Unit circle
- S^2 , *see* Sphere
- S^3 , *see* 3-sphere, *see* 3-sphere
- S_4 , 119
- S_5 , 108
- S_n , *see* Symmetric group
- \mathbf{C} , *see* Complex numbers
- $\mathbf{C}[x]$, *see* Complex polynomials
- \mathbf{N} , *see* Natural numbers
- Φ_p , *see* Prime cyclotomic polynomials
- \mathbf{Q} , *see* Rational numbers
- \mathbf{R} , *see* Real numbers
- $\mathbf{R}[x]$, *see* Real polynomials
- \Rightarrow , *see* Implication
- \mathbf{Z} , *see* Integers
- $\mathbf{Z}[i]$, *see* Gaussian integers
- \mathbf{Z}_3P^2 , 128
- \mathbf{Z}_n , *see* Integers modulo n
- \triangleright^{-1} back through, 117
- \therefore , since, 143
- \cong , *see* Isomorphism
- \mid , *see* Divisibility
- \equiv , *see* Equivalence relation
- $\equiv \pmod{n}$, *see* Congruence modulo n
- \exists , *see* Existential quantification
- \forall , *see* Universal quantification
- \iff , *see* Bi-implication
- \wedge , *see* Conjunction
- \vee , *see* Disjunction
- \mathbf{H} , *see* Quaternions
- \mathcal{C} , *see* Category
- \mathcal{G} , *see* Category of groups
- \mathcal{R} , *see* Category of rings
- \mathcal{S} , *see* Category of sets
- \neg , *see* Negation (logical)
- \oplus , *see* Direct sum
- \mathcal{P} , *see* Powerset
- \preceq , *see* Partial order
- \therefore , therefore, 143
- \triangleright through, 117
- p -group, 138
- p -primary component, 138
- p -primary group, 138
- 3-Sphere, 54
- Abel, Neils Henrik (1802–1829), 134
- Abelian group, 7, 99
 - finite, 137–140
- ACC (ascending chain condition), 83
- Affine geometry, 127
- al-Khwārizmī (ca. 780–ca. 850), 1
- Algebra, 1
 - Boolean, 157
 - Cayley, 52
 - division, 52
- Algebraic field extension, 42
- Algebraic fields, 42
- Algebraic integer, 42
- Algebraic number, 42
- Algebraic structure, 2–11
- Algebraically closed field, 89
- Algorithm
 - Brahmagupta’s, 62
 - division, 84
 - Euclidean, 19–20, 86
 - extended Euclidean, 22
 - Qin Jiushao’s, 62
- Alternating group A_n , 107, 112
- Antiautomorphism, 52
- Antisymmetry, 66
- Archimedean ordered field, 47–48
- Archimedes of Syracuse (ca. 287–212 B.C.E.), 47

- Aristotle (384–322 B.C.E.), 17
- Arrow, *see* morphism
- Ascending chain condition, 83
- Associativity, 2
- Automorphism, 14
field, 43
- Axiom
of choice, 153, 158–160
of extensionality, 151
of infinity, 152
of pairing, 151
of powersets, 152
of regularity, 152
of replacement, 152
of separation, 151
of union, 152
- Axioms
Dedekind/Peano, 15
field, 31
group, 99
of set theory, 151–153
ring, 55
- Back through \triangleright^{-1} , 117
- Bernays, Paul Isaak (1888–1977), 153
- Bernstein, Felix (1878–1956), 148
- Bi-implication \iff , 144
- Bijection, 150
- Binary operation, 2
- Binary order relation, 46
- Bombelli, Rafael (1526–1572), 87
- Boole, George (1815–1864), 56, 63
- Boolean algebra, 157
- Boolean ring, 56, 63–67
- Bound
greatest lower, 154
least upper, 154
lower, 154
upper, 154
- Brahmagupta (598–670), 62
- Brahmagupta’s algorithm, 62
- Cancellation, 58
- Canonical function, 151
- Canonical homomorphism, 38
- Cardano, Gerolamo (1501–1576), 87
- Cardinality, 101, 148
- Category, 69–74
coproduct, 136
final object, 72, 115
generic, \mathcal{C} , 13
initial object, 72, 115
of Abelian groups \mathcal{A} , 136–137
of fields, 71
of groups \mathcal{G} , 71, 115
of rings \mathcal{R} , 71–74
of sets \mathcal{S} , 71
- Cauchy sequence, 49
- Cauchy, Augustin Louis (1789–1857), 49
- Cayley algebra, 52
- Cayley’s theorem, 110–112
- Cayley, Arthur (1821–1895), 1, 5, 52, 110
- Center of a group, 101
- Centralizer, 101
- Chain, 158
- Characteristic
of a field, 41
of a ring, 38
of an integral domain, 58, 69
- Charles Hermite (1882–1901), 42
- Chinese remainder theorem, 60–63, 77, 102
- Circle
unit, 9, 89
- Codomain, 70, 149
- Cohen, Paul (1934–2007), 158
- Commutative diagram, 70
- Commutative group, *see* Abelian group
- Commutative ring, 55
- Commutativity, 2
- Commutator
of two group elements, 101
subgroup, 101
- Comparable elements, 154
- Complete ordered field, 49–50
- Complex conjugation, 14, 43, 89
- Complex numbers \mathbf{C} , 2, 4, 43–44, 87–89
- Complex polynomials $\mathbf{C}[x]$, 87
- Composite number, 17
- Composition
of functions, 149
of homomorphisms, 13
of morphisms, 70
- Composition factor, 131
- Composition series, 131
- Congruence
group, 121
ring, 76
- Congruence class, 37, 76, 121

- Congruence modulo n , 5, 36
- Conjecture, 143
- Conjugacy class, 116
- Conjugate
 - element in a group, 115
 - subgroup, 116
- Conjugation
 - complex, 14, 43, 89
 - for a quadratic extension field, 43
 - quaternion, 51
- Conjunction \wedge , 144
- Content of a polynomial, 93
- Contraction, 125
- Conway, John H., 47
- Coprime, 17
- Coproduct
 - in a category, 136
- Core of a group, 117
- Corollary, 143
- Correspondence theorem for groups, 123
- Coset, 103–104
- Countable set, 149
- Cross product, 53
- CRT, *see* Chinese remainder theorem
- Cubic equation, 87
- Cubic polynomial, 91
- Cycle notation, 105
- Cyclic field, 39–40
- Cyclic group, 101–102
- Cyclotomic polynomial, 29

- d’Alembert, Jean le Rond (1717–1783), 89
- Dave’s Short Course on Complex Numbers, 4
- de Foncenex, François Daviet (1734–1799), 89
- De Morgan’s laws, 147
- De Morgan, Augustus (1806–1871), 147
- Dedekind cut, 48
- Dedekind, Richard (1831–1916), 15
- Dedekind/Peano axioms, 15
- Desargues
 - Girard (1591–1661), 129
- Determinants
 - as group homomorphisms, 125
- Diagram
 - commutative, 70
- Dihedral group D_n , 108, 111
- Dilation, 125
- Direct sum \oplus , 102
- Disjoint
 - pairwise, 146
 - sets, 146
 - union, 137
- Disjunction \vee , 144
- Distributivity, 3, 33, 67
- Divisibility \cong , 80–81
- Divisibility $|$, 16–17
- Division algorithm, 84
 - for polynomials, 26
- Division ring, 10, 50–54
- Domain, 70, 149
 - Euclidean, 84–87
 - integral, 57–60, 78, 80
 - principal ideal, 82–84, 86
 - unique factorization, 81–82, 84
- Dot product, 53, *see* inner product
- Dyadic rational, 69

- ED, *see* Euclidean domain
- Eilenberg, Samuel (1913–1998), 69
- Eisenstein integers, 59, 86
- Eisenstein’s criterion, 92–95
- Eisenstein, Gotthold (1823–1852), 59, 86, 92
- Element
 - identity, 3
 - initial, 15
 - inverse, 3
 - irreducible, 80–82
 - maximal, 158
 - order of, 101
 - positive and negative, 45
 - prime, 81, 82
- Elements, 145
- Elements of Euclid, 16–18
- Endomorphism, 14
- Epimorphism, 13, 73
- Equivalence class, 150–151
- Equivalence relation \equiv , 67, 77, 150–151
- Euclid of Alexandria (fl. ca. 300 B.C.E.), 16–19, 47
- Euclidean algorithm, 19–20, 86
- Euclidean domain, 84–87
- Euclidean geometry, 117, 127
- Euclidean valuation, 84
- Eudoxus (fl. 350 B.C.E.), 47
- Euler’s circle group, 9
- Euler’s identity, 89
- Euler, Leonhard (1707–1783), 9, 19, 43, 89, 146
- Even permutation, 106–107

- Existential quantification \exists , 144
 Expansion, 125
 Extended Euclidean algorithm, 22
 Extension field, 41–43

 Factor theorem, 27
 Fano plane, 128
 Fano, Gino (1871–1952), 128
 Fermat, Pierre de (1607–1665), 16
 Field, 2, 4, 31–54, 79
 - algebraic, 42
 - algebraic number, 42
 - algebraically closed, 89
 - Archimedean, 47–48
 - axioms, 31
 - category, 71
 - complete ordered, 49–50
 - definition, 4, 31
 - extension, 41–43
 - homomorphism, 14
 - isomorphism, 12
 - number, 42, 97–98
 - of complex numbers, 43–44
 - of rational functions, 35, 69
 - of rational numbers, 34, 67–69
 - ordered, 45–50
 - prime, 39–41
 - skew, 10, 50–54
 Field extension
 - algebraic, 42
 - quadratic, 41–45, 77
 - transcendental, 42
 Final object, 72, 115
 Finding one, 63
 Finite Abelian group, 137–140
 Finite group, 103, 104, 112–115
 First isomorphism theorem for groups, 122
 First isomorphism theorem for rings, 78
 Fixed point, 104, 105
 Four squares identity, 52
 Fraenkel, Abraham (1891–1965), 151
 Free Boolean ring, 65
 Free group, 115
 Frobenius endomorphism, 41
 Frobenius, Ferdinand Georg (1849–1917), 41, 52, 133
 FTA, *see* Fundamental theorem of algebra
 Function, 149
 - canonical, 151
 - choice, 158
 - codomain, 149
 - composition, 149
 - domain, 149
 - graph, 149
 - identity, 13, 149
 - image, 149
 - inclusion, 149
 - injective, 13, 73, 149
 - inverse, 150
 - preimage, 149
 - projection, 151
 - rational, 35, 69
 - successor, 15
 - surjective, 13, 73, 149
 Fundamental theorem
 - of algebra, 88–89
 - of arithmetic, 22–24
 - of finite Abelian groups, 140
 Gödel, Kurt (1906–1978), 153, 158
 Galois field $GF(p^n)$, 40, 129
 Galois, Evariste (1811–1832), 40, 133
 Gauss’s lemma, 92–95
 Gauss, Carl Friedrich (1777–1855), 5, 59, 69, 85, 89, 93
 Gaussian integers, 59
 Gaussian integers $\mathbf{Z}[i]$, 69, 85
 GCD, *see* greatest common divisor
 General linear group $GL_n(R)$ $GL_n(R)$, 7, 125–126
 Geodesic, 118
 Geodesics, 117
 Geometry
 - affine, 127
 - Euclidean, 117, 127
 Georg Cantor (1845–1918), 149
 Girard, Albert (1595–1632), 88
 Gorenstein, Daniel (1923–1992), 112
 Grassmann, Hermann (1809–1977), 1
 Graves, John T. (1806–1870), 52
 Greatest common divisor, 19, 22, 80
 Greatest lower bound, 154
 Group, 2, 6–10, 99–140
 - Abelian, 7, 99, 134–140
 - alternating, 107, 112
 - axioms, 99
 - category, 71, 115
 - center, 101

- circle, 9
- core, 117
- cyclic, 7, 101–102, 111
- definition, 6
- dihedral, 108, 111
- finite, 103, 104, 112–115
- finite Abelian, 137–140
- free, 115
- Frobenius, 133
- general linear, 7, 125–126
- homomorphism, 13
- isomorphism, 12
- Klein 4-group, 113
- linear, 124–130
- of units in a ring, 7
- order, 7, 101, 103
- orthogonal, 126
- presentation, 110
- primary, 138
- projective linear, 127, 130
- projective special linear, 130
- quaternion, 113, 119
- quotient, 121–123
- simple, 131–134
- solvable, 133
- special linear, 126
- sporadic, 133
- symmetric, 104–107, 110, 112, 116, 119
- unitary, 127
- Group action
 - transitive, 133
- Group ring, 58
- Hölder, Otto (1859–1937), 131
- Hamilton, William Rowan (1805–1865), 1, 50
- Hasse diagram, 16, 154
- Hasse, Helmut (1898–1979), 16, 154
- Heptahedron, 134
- Hermite, Charles (1822–1901), 127
- Hermitian, 127
- Hom set, 70
- Homomorphism, 12–13
 - field, 14
 - group, 13
 - ring, 13, 71
- Hyperbolic space, 118
- Ideal, 74–79
 - generated by a set, 75
 - maximal, 78–79
 - prime, 78–79, 82
 - principal, 75
 - proper, 75
 - trivial, 75
- Idempotent element, 64
- Identity element, 3
- Identity morphism, 70
- Image, 149
- Implication \Rightarrow , 144
- Inclusion, 14
- Inclusion function, 149
- Index of a subgroup, 103
- Infimum, *see* Greatest lower bound
- Initial element, 15
- Initial object, 72, 115
- Injection, 13, 73, 149
- Inner product, 53, 126
- Integer
 - algebraic, 42
- Integers
 - Eisenstein, 59, 86
 - Gaussian, 59, 69, 85
- Integers \mathbf{Z} , 2, 4, 5, 69, 72
- Integers modulo n , \mathbf{Z}_n , 5, 7, 35–39, 56, 60–63, 74, 77, 102
 - definition, 36
- Integral domain, 57–60, 78, 80
- Internal direct product, 123–124
- Internal direct sum, 136
- Intersection, 146, 147
 - of subgroups, 100
- Inverse element, 3
- Inverse function, 150
- Inversive space, 118
- Invertible element, *see* unit
- Involution, 101
- Involutory quandle, 117
- Irreducibility test
 - Eisenstein’s criterion, 94
 - modulo p , 93
- Irreducible element, 80–82
- Irreducible polynomial, 87, 89–97
- Isomorphism \cong , 11–12, 71, 74
 - field, 12
 - group, 12
 - ring, 11
- Isomorphism theorem
 - first for groups, 122

- first for rings, 78
 - second for groups, 123
 - third for groups, 123
- Join, 66
- Jordan, Camille (1838–1922), 122, 131
- Jordan-Hölder theorem, 131–133
- Joyce, David, 4, 16, 145
- Kelland, Philip (1808–1879), 50
- Kernel
 - of a group homomorphism, 120–124
 - of a ring homomorphism, 74
- Klein, Felix (1849–1925), 113
- Knuth, Donald, 47
- Krull's theorem, 79
- Krull, Wolfgang (1899–1971), 79
- Lagrange's theorem, 103–104
- Lagrange, Joseph-Louis (1736–1813), 52, 89, 103
- Laplace, Pierre-Simon (1749–1827), 89
- Latin square, 9
- Lattice, 67, 155
 - distributive, 67, 156
 - modular, 156
- Least common denominator, 22
- Least common multiple, 22
- Least upper bound, 154
- Leibniz, Gottfried Wilhelm (1646–1716), 89, 146
- Lemma, 143
- Lexicographic ordering, 159
- Lindemann, Ferdinand von (1852–1939), 42
- Linear group, 124–130
- Linear order, 153
- Linear transformation, 7, 44, 53, 124–127
- Llull, Ramon (ca. 1232–ca. 1315), 146
- Localization, 69
- Logical symbols, 143
- Loos, Ottmar, 117
- Lower bound, 154
- Mac Lane, Saunders (1909–2005), 69
- Map, 12, *see* morphism, 149
- Mathematical induction, 15
 - strong form, 23
- Mathieu, Émile Léonard (1835–1890), 133
- Matrix
 - unimodular, 126
 - unitary, 127
- Matrix representation
 - of \mathbf{C} , 43
 - of \mathbf{H} , 53
- Matrix ring, 5, 35, 43, 124–130
- Maximal ideal, 78–79
- Meet, 66
- Membership, 145
- Minimization principle, 15
- Module, 11
- Modulo p irreducibility test, 93
- Monomorphism, 13, 73
- Morphism, 12, 70
- Moufang, Ruth (1905–1977), 52
- Multiplicative function, 19
- Multiplicative group of units, 7
- Natural numbers \mathbf{N} , 2, 15
- Negation (logical) \neg , 144
- Neutral element, *see* identity element
- Noether, Emmy (1882–1935), 1
- Noether, Emmy Amalie (1882–1935), 83
- Noetherian ring, 83
- Norm
 - of a complex number, 43
 - of a quaternion, 51
- Normal subgroup, 120–124
- Number
 - algebraic, 42
 - complex, 2, 4, 43–44, 87–89
 - composite, 17
 - greatest common divisor, 22
 - integer, 2, 4, 5, 69, 72
 - natural, 2, 15, 144
 - prime, 17–19, 22–39
 - rational, 2, 4, 34, 67–69
 - real, 2, 45–50, 89–90, 145
 - relatively prime, 17, 19–20, 24, 39, 61, 102, 137
 - surreal, 47
 - transcendental, 42
 - whole, *see* integers
- Number field, 42, 97–98
- Number theory, 15–25
- Object, 70
- Octonions, 52
- Odd permutation, 106–107
- One-to-one correspondence, *see* bijection
- One-to-one function, *see* injection
- Onto function, *see* surjection

- Operation, 2–3
 - associative, 2
 - binary, 2
 - commutative, 2
 - unary, 2
- Order
 - lexicographic, 159
 - linear, 153
 - of a group, 7, 101, 103
 - of a prime in a number, 25
 - of an element in a group, 101
 - partial, 66, 153–155
 - total, 153
- Ordered field, 45–50
 - Archimedean, 47–48
 - complete, 49–50
- Orthogonal group, 126
- Orthogonal transformation, 126
- Outer product, 53

- Pairwise relatively prime numbers, 22
- Pappus of Alexandria (ca. 290–ca. 350), 129
- Partial order \preceq , 66, 153–155
- Partition, 151
- Partition of a number, 138
- Partition of a set, 148, 150–151
- Peano, Giuseppe (1858–1932), 15
- Permutation, 104
 - even and odd, 106–107
- Philolaus (470–385 B.C.E.), 17
- PID, *see* principal ideal domain
- Polynomial, 25–29
 - complex, 87
 - content, 93
 - cubic, 91
 - cyclotomic, 29
 - irreducible, 87, 89–97
 - monic, 25
 - prime cyclotomic, 95
 - primitive, 92
 - quadratic, 90
 - rational root theorem, 91
 - real, 89
 - root, 26
- Polynomial evaluation, 13, 73
- Polynomial ring, 5, 26, 73, 85–97
- Poset, 153–155
- Powerset \mathcal{O} , 56, 64, 148
- Pre-order, 154
- Preimage, 149
- Presentation by generators and relations, 110
- Primary component, 138
- Primary decomposition theorem, 138
- Primary group, 138
- Prime cyclotomic polynomials Φ_p , 95
- Prime element, 81, 82
- Prime field, 39–41
- Prime ideal, 78–79, 82
- Prime number, 17–19, 22–39
 - infinitely many, 18
- Primitive polynomial, 92
- Primitive root of unity, 28
- Primitive roots of unity, 95
- Principal ideal, 75
- Principal ideal domain, 82–84, 86
- Principle of infinite descent, 16
- Product
 - in a category, 71
 - internal direct, 123–124
 - of groups, 102
 - of rings, 57, 71
 - of sets, 148
 - semidirect, 133
- Products of subsets in a group, 104
- Projection, 38, 151
- Projective linear group $PGL_n(F)$, 127, 130
- Projective plane
 - Desarguesian, 129
 - finite, 129
 - Pappian, 129
- Projective space, 118, 127
- Projective special linear group $PSL_n(F)$, 130

- Q.E.D., 143
- Qin Jiushao (1202–1261), 62
- Qin Jiushao’s algorithm, 62
- Quadratic field extension, 41–45, 77
- Quadratic polynomial, 90
- Quandle, 10, 11, 117
 - involutory, 117
 - with geodesics, 117
- Quaternion group, 113, 119
- Quaternions \mathbf{H} , 10, 50–54
 - unit, 54
- Quotient group, 121–123
- Quotient ring, 76–79
- Quotient set, 37, 68, 76, 151

- Radian, 89
- Rational function, 35, 69
- Rational numbers, 2, 4, 34, 67–69
- Rational root theorem, 91
- Real numbers, 45–50, 89–90
- Real numbers \mathbf{R} , 2
- Real polynomials $\mathbf{R}[x]$, 89
- Reducible, 80
- Reflection, 125
- Reflexivity, 66, 150
- Relation, 150
 - antisymmetric, 66
 - binary order, 46
 - equivalence, 67, 77, 150–151
 - partial order, *see* Partial order
 - reflexive, 66, 150
 - symmetric, 150
 - transitive, 16, 66, 150
- Relatively prime, 17, 19–20, 24, 39, 61, 102, 137
 - pairwise, 22
- Remainder theorem, 27
- Residue, 6
- Ring, 2, 5–6, 55–98
 - algebraic integers, 97–98
 - axioms, 55
 - Boolean, 56, 63–67
 - category, 71–74
 - commutative, 55
 - cyclic, 35–38, 56, 77
 - definition, 5
 - division, 10, 50–54
 - free Boolean, 65
 - homomorphism, 13, 71
 - isomorphism, 11
 - matrix, 5, 35, 124–130
 - Noetherian, 83
 - of integers, *see* integers
 - of polynomials, 5, 26, 73, 85–97
 - quotient, 76–79
 - trivial, 72
- Root of unity, 28–29, 95
 - primitive, 28, 95
- Rotation, 125

- Scalar, 53
- Scalar product, 53
- Schröder, Ernst (1841–1902), 148
- Second isomorphism theorem for groups, 123
- Semidirect product, 133

- Set, 15, 144–149
 - category 71
 - complement, 147
 - countable, 149
 - difference, 147
 - element, 145
 - finite, 11
 - infinite, 15
 - intersection, 146, 147
 - membership, 145
 - operation on, 2–3
 - partially ordered, 153–155
 - partition, 148, 151
 - permutation, 104
 - power, 56, 64, 148
 - product of, 148
 - quotient, 37, 68, 76, 151
 - singleton, 152
 - subset, 145
 - uncountable, 149
 - underlying, 3, 31, 55, 99
 - union, 146, 147
- Set theory
 - axioms, 151–153
- Shear, 126
- Simple group, 131–134
- Simply infinite, 15
- Singleton set, 152
- Skew field, 10, 50–54
- Solvable group, 133
- Space
 - hyperbolic, 118
 - inversive, 118
 - projective, 118, 127
- Special linear group, 126
- Sphere, 118
- Sphere S^2 , 10
- Structure
 - algebraic, 2–11
- Subfield, 34
 - definition, 34
- Subgroup, 14, 100–104
 - commutator, 101
 - conjugate, 116
 - generated by a set, 101
 - generated by an element, 101
 - index, 103
 - normal, 120–124
 - of \mathbf{Z} , 101

- of S_3 , 107
 - proper, 100
 - trivial, 100
- Subring, 14, 58
- Subset, 145
- Substitution
 - Tschirnhaus, 91
- Successor function, 15
- Sun Zi (fl. 400), 62
- Supremum, *see* Least upper bound
- Surjection, 13, 73, 149
- Surreal numbers, 47
- Sylvester, James Joseph (1814–1897), 5
- Symmetric group S_n , 104–107, 110, 112, 116, 119
- Symmetric space, 117
- Symmetries
 - of a cube, 108
 - of a pentagon, 108, 111, 119
 - of a tetrahedron, 108
 - of a triangle, 8
- Symmetry, 150

- Tait, Peter Guthrie (1831–1901), 50
- Tartaglia, Nicolo Fontana (1500–1557), 87
- Theorem, 143
- Third isomorphism theorem for groups, 123
- Through \triangleright , 117
- Thymaridas (400–350 B.C.E.), 17
- Torus, 134
- Total order, 153
- Totient function, 19, 29, 39
- Transcendence
 - of π , 42
 - of e , 42
- Transcendental field extensions, 42
- Transcendental number, 42
- Transformation
 - linear, 7, 44, 53, 124–127
- Transitive group action, 133
- Transitivity, 16, 46, 66, 150
- Transposition, 105, 106
- Trichotomy, 46
- Trivial ring, 72
- Tschirnhaus substitution, 91
- Tschirnhaus, Ehrenfried Walther von (1651–1708), 91

- UFD, *see* unique factorization domain
- Unary operation, 2

- Uncountable set, 149
- Underlying set, 3
- Unimodular matrix, 126
- Union, 146, 147
 - disjoint, 137
- Unique factorization domain, 81–82, 84, 95
- Unique factorization theorem, 22–24
- Unit
 - circle, 9
 - in \mathbf{Z}_n , 19
 - in a ring, 7
- Unit circle S^1 , 89
- Unitary group, 127
- Unitary matrix, 127
- Unitary transformation, 127
- Unity
 - root of, 28–29, 95
- Universal property
 - of an infinite cyclic group, 115
 - of coproducts, 136
 - of final objects, 72, 115
 - of free groups, 115
 - of initial objects, 72, 115
 - of products, 71
 - of the ring \mathbf{Z} , 72, 73
- Universal quantification \forall , 144
- Upper bound, 154

- Valuation
 - Euclidean, 84
- Vector, 53
- Vector product, 53
- Vector space, 35
- Venn diagram, 146
- Venn, John (1834–1923), 146
- Viète, François (1540–1603), 91
- von Neumann, John (1903–1957), 153

- Waring, Edward (1736–1798), 50
- Wedderburn, Joseph (1882–1948), 59
- Weierstrass, Karl (1815–1897), 42
- Well-ordering principle, 15, 159

- Zermelo, Ernst (1871–1953), 151
- Zero-divisor, 58
- Zorn's lemma, 79, 158–160
- Zorn, Max August (1906–1993), 158