

Finite Fields and Pseudo-Random Number Generation

Carl Offner

February 4, 2007

Contents

1	Algebraic Preliminaries	3
1.1	Arithmetic	3
1.1.1	Division	4
1.1.2	Primes	4
1.1.3	Congruence	6
1.2	Binomial coefficients	8
1.3	Groups	12
1.3.1	Definition and examples	12
1.3.2	Subgroups	15
1.3.3	Products of groups	17
1.3.4	The order of an element of a group	19
1.3.5	Orders of elements in abelian groups	21
1.4	Rings	23
1.5	Fields	24
1.6	Vector spaces	26
1.6.1	Vector spaces and dimensionality	26
1.6.2	Linear functions	30
1.6.3	Dot products and adjoint transformations	32
1.6.4	An application of vector spaces	33
1.7	Polynomials over a field	36
1.8	Factoring polynomials	38
1.8.1	Division of polynomials	38
1.8.2	Linear factors of polynomials	40

1.8.3	Formal derivatives	43
2	Finite Fields	45
2.1	Irreducible polynomials and minimal polynomials	45
2.2	Algebraic field extensions	46
2.3	Finite fields	54
2.4	Primitive polynomials	58
3	Some Random Number Generators	61
3.1	Linear pseudo-congruential generators	61
3.2	Higher order linear recursive generators	62
3.3	Inversive congruential generators	66

Introduction

The purpose of this set of notes is to show, as simply as possible, how the theory of finite fields applies to certain commonly used pseudo-random number generators. Only those parts of the theory of finite fields that are needed for this purpose are presented, and the development of the algebraic theory needed for this is greatly simplified for this purpose. I have tried to compose these notes in a sympathetic manner, so as to be readable to someone who has not read mathematics in a while.

The notes only go so far as to show how finite fields enter into the theory of pseudo-random number generators. I do not pursue the crucial and interesting investigations into the quality and efficiency of such generators.

While all of the algebraic material here is quite standard, the following points may be noted:

- Undoubtedly the most difficult part of this subject is the construction of field extensions in Chapter 2. This is not because the construction itself is complicated—it is really quite simple. The reason is that the arguments needed to justify the correctness of the construction involve isomorphisms of quotient constructions (to use the technically correct terms); such arguments tend to be confusing to beginners. I have tried to make this section quite discursive and have given a number of examples, taken verbatim from Lidl and Niederreiter.
- Since I am not looking to present things in much generality, but am only aiming at the theory of finite fields, I completely eliminated the standard but sophisticated arguments involving ideals in rings. These come up in two contexts:
 1. To get unique prime factorization. Even in general, ideals are not needed for this, although their use undoubtedly makes the theory more elegant. The traditional proof as given in Hardy and Wright, however, works just fine here. I didn't bother writing it out.
 2. To show that adjoining the root of an irreducible polynomial really does yield a field. In modern treatments this is usually approached

by way of maximal ideals. Now actually, there is a simple traditional proof that follows directly from the division algorithm, so ideals are not really needed here either. In the case of finite fields, however, since everything is finite (and not just finite-dimensional), we can avoid even this traditional argument and reduce the whole matter to the kind of counting argument that is used over and over in these notes.

- In a similar way, I have given a simplified proof of the fact that any two finite fields having the same number of elements are isomorphic. The proof of this result is usually made to rest on the uniqueness of splitting fields. Again in this case, I was able to avoid using this fact.

None of these simplifications is at all original. But practically all algebra texts prove things in more generality and so cannot use these simple arguments.

- The proof in Chapter 3 that primitive polynomials yield higher-order recursive random number generators with maximal period is greatly simplified from the two references I have seen.

The reader is assumed to know (or at least to be willing to remember) some basic linear algebra. What is required is reviewed in the section on vector spaces in Chapter 1. Other than that, these notes are pretty much self-contained.

Chapter 1

Algebraic Preliminaries

There are many mathematical structures that consist of a set with one or more binary operations that seem similar to ordinary addition and/or multiplication. Mathematicians have abstracted these structures so that their common properties can be handled all at once. The four most important of these structures are groups, rings, fields, and vector spaces¹. In this chapter we give brief introductions to these four structures. First, however, we review some arithmetic.

1.1 Arithmetic

Here is some standard notation that we often use:

\mathbf{N} denotes the set of positive integers (i.e., the integers *greater than 0*). The letter \mathbf{N} comes from the word “natural”—the set \mathbf{N} is sometimes referred to as the set of natural numbers; i.e., the set of counting numbers that everyone learns in first grade.

\mathbf{Z} denotes the set of all integers (including 0 and the negative integers). The letter \mathbf{Z} comes from the German word “Zahlen”, which means “numbers”.

\mathbf{R} denotes the set of all real numbers.

\mathbf{C} denotes the set of all complex numbers.

The notation $b|a$ means “ b divides a ” (or, as my students would have put it, “ b guzzinta a ”). For example, $6|18$.

¹The terms “group”, “ring”, and “field” were not picked because of any intuitive meaning. In particular, there is nothing ring-like about a ring.

1.1.1 Division

When we first learn division, we learn to say “7 divided by 3 is 2 with a remainder of 1.” This is the same as writing

$$7 = 2 \cdot 3 + 1$$

The remainder is always ≥ 0 and strictly less than the number we are dividing by (which in this case is 3).

In general, if a and b are two integers and $b > 0$, we can always divide a by b to get a quotient q and a remainder r . We always determine r so that $0 \leq r < b$, and we have

$$a = qb + r$$

Thus, if we are dividing by 7, we have

$$14 = 2 \cdot 7 + 0$$

$$9 = 1 \cdot 7 + 2$$

$$2 = 0 \cdot 7 + 2$$

$$-5 = -1 \cdot 7 + 2$$

Note that $b|a$ (b divides a) if and only if $r = 0$.

1.1.2 Primes

The prime numbers in \mathbf{N} have the following two characteristic properties:

- A prime p cannot be factored into two non-unit factors (i.e., into two factors neither of which is 1).
- If p is a prime, and if a and b are any two numbers in \mathbf{N} , and if p divides the product ab (remember that the notation for this is $p|ab$), then either $p|a$ or $p|b$ (or both, of course).

The first property is conventionally regarded as a defining property for primes (with the additional statement that 1 is not a prime). But the second property is equally characteristic of primes, and is absolutely key. For instance, from the fact that 3 divides a product ab we can conclude that either 3 divides a or 3 divides b . On the other hand, 6 (which is not a prime) divides the product $14 \cdot 27 = 378$, but 6 does not divide either 14 or 27.

The *Fundamental Theorem of Arithmetic* states that every positive integer is the unique product of primes. That is given any positive integer a , we can write

$$a = p_1^{e_1} p_2^{e_2} \cdots p_n^{e_n}$$

where the exponents $\{e_1, e_2, \dots, e_n\}$ are determined uniquely by a . For example,

$$2352 = 2^4 \cdot 3^1 \cdot 7^2$$

We could actually write this as a product over all the primes, by making the remainder of the exponents all 0:

$$2352 = 2^4 \cdot 3^1 \cdot 5^0 \cdot 7^2 \cdot 11^0 \cdot 13^0 \cdot 17^0 \cdots$$

This is the usual way that mathematicians write this, except that we write it using the product notation:

$$a = \prod_{p_i \text{ prime}} p_i^{e_i}$$

The reason we agree that 1 is not a prime is to allow for the statement of uniqueness in the Fundamental Theorem. If 1 were a prime, we could always include an arbitrary number of factors of 1 into the product, and so the product representation would not be unique.

The Fundamental Theorem illuminates many problems of divisibility. For example, the reason why $6 \mid 14 \cdot 27$ even though 6 does not divide either 14 or 27 is that one of the prime factors (2) of 6 divides 14, and the other one (3) divides 27, but neither 14 nor 27 is divisible by both.

In general, if

$$\begin{aligned} a &= \prod_{p_i \text{ prime}} p_i^{\alpha_i} \\ b &= \prod_{p_i \text{ prime}} p_i^{\beta_i} \end{aligned}$$

then a divides b if and only if $\alpha_i \leq \beta_i$ for all i .

It follows that, given two numbers a and b , their greatest common factor $\text{gcf}(a, b)$ and their least common multiple $\text{lcm}(a, b)$ can be computed as follows:

$$\begin{aligned} \text{gcf}(a, b) &= \prod_{p_i \text{ prime}} p_i^{\min(\alpha_i, \beta_i)} \\ \text{lcm}(a, b) &= \prod_{p_i \text{ prime}} p_i^{\max(\alpha_i, \beta_i)} \end{aligned}$$

For instance, the least common multiple of

$$2^3 \cdot 5^4 \cdot 11 \quad \text{and} \quad 2^2 \cdot 5^5 \cdot 7^3$$

is

$$2^3 \cdot 5^5 \cdot 7^3 \cdot 11$$

These formulas are not often useful for computation, because usually we do not know the prime factorization of a number. Usually, by far the best way of computing the greatest common factor is by use of the Euclidean algorithm. Nevertheless, these formulas are important in understanding and proving things.

For another example, to say that a and b are relatively prime is just to say that $\text{gcf}(a, b) = 1$; i.e., no prime divides both a and b . Thus, if a and b are relatively prime and n is some number such that $a|nb$, then by looking at the factorization into primes, we see that a must divide n .

1.1.3 Congruence

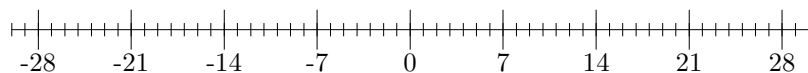
The notation

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

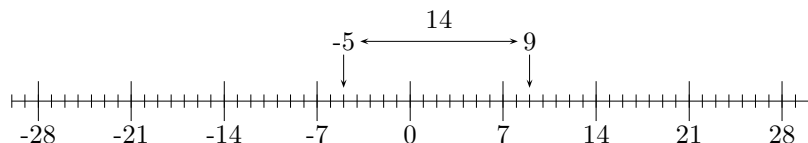
simply means that $m|a - b$. Another way to think of this is that it means that a and b each have the same remainder when divided by m . The notation is read “ a is congruent to b mod m ”, which is short for “ a is congruent to b with respect to the modulus m ”. This notation is due to Gauss.

In particular, $a \equiv 0 \pmod{m}$ is just another way of saying that $m|a$.

A good way to think of congruence is by looking at the number line. Suppose we want to consider congruence mod 7. We draw a number line, and label the multiples of 7:



Then 9 and -5 are congruent (mod 7) (in symbols, $9 \equiv -5 \pmod{7}$) because the distance between them on the number line is a multiple of 7:



Congruence is an equivalence relation. Numbers that are congruent are said to be in the same congruence class. So for instance the congruence class containing

the number 9 is the set

$$\{\dots, -33, -26, -19, -12, -5, 2, 9, 16, 23, 30, 37, \dots\}$$

There are 7 congruence classes (mod 7). Each congruence class contains one number in the set $\{0, 1, 2, 3, 4, 5, 6\}$. We often identify these numbers with the congruence classes (mod 7).

The phrase “the integers mod 7” refers to these 7 congruence classes.

Here are some important properties of congruences. We will write them as theorems so that we can refer to them later, although actually they are pretty elementary:

1.1 Theorem *If*

$$a_1 \equiv b_1 \pmod{m}$$

and

$$a_2 \equiv b_2 \pmod{m}$$

then

$$a_1 + a_2 \equiv b_1 + b_2 \pmod{m}$$

This is pretty obvious if you look at the number line. Another way to see it is as follows: we know by assumption that $m|a_1 - b_1$ and $m|a_2 - b_2$. Therefore, m divides their sum: $m|(a_1 - b_1) + (a_2 - b_2)$, which is the same as saying $m|(a_1 + a_2) - (b_1 + b_2)$, and this in turn is just the statement that

$$a_1 + a_2 \equiv b_1 + b_2 \pmod{m}$$

1.2 Theorem *If*

$$a_1 \equiv b_1 \pmod{m}$$

and

$$a_2 \equiv b_2 \pmod{m}$$

then

$$a_1 a_2 \equiv b_1 b_2 \pmod{m}$$

This is true because the first statement says that $a_1 = b_1 + cm$ for some integer c . (That is, a_1 and b_1 differ by a multiple of m .) Similarly, the second statement says that $a_2 = b_2 + dm$ for some integer d . Thus, we have

$$a_1 a_2 = (b_1 + cm)(b_2 + dm) = b_1 b_2 + m(cb_2 + db_1 + cdm)$$

That is, a_1a_2 and b_1b_2 differ by a multiple of m , which is the conclusion we needed to reach.

1.3 Theorem *As a simple special case of the previous property, if*

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

then

$$an \equiv bn \pmod{m}$$

This just says that if $m|a - b$ then $m|(a - b)n$.

The converse is not true in general: for instance, $5 \cdot 3 \equiv 7 \cdot 3 \pmod{6}$, but $5 \not\equiv 7 \pmod{6}$. However, there is an important case in which the converse holds:

1.4 Theorem *If*

$$an \equiv bn \pmod{m}$$

and if m and n are relatively prime, then

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

This is because, as we have seen, if $m|(a - b)n$ and m and n are relatively prime, then $m|(a - b)$.

1.2 Binomial coefficients

In the familiar Pascal's triangle

$$\begin{array}{cccccc}
 & & & & 1 & & & & & & \\
 & & & & & 1 & & 1 & & & \\
 & & & & & & 1 & & 2 & & 1 & \\
 & & & & & & & 1 & & 3 & & 3 & & 1 & \\
 & & & & & & & & 1 & & 4 & & 6 & & 4 & & 1 & \\
 & & & & & & & & & 1 & & 5 & & 10 & & 10 & & 5 & & 1 & \\
 & & & & & & & & & & & & & \vdots & & & & & & &
 \end{array}$$

each number is the sum of the two numbers diagonally above it, and the end numbers on each row are defined to be 1. The rows are numbered starting with 0, and the elements in each row are also numbered starting with 0. The r^{th} element in the n^{th} row is denoted

$$\binom{n}{r}$$

in the expansion of the binomial power $(a + b)^n$:

$$(a + b)^0 = 1a^0b^0$$

$$(a + b)^1 = 1a^1b^0 + 1a^0b^1$$

$$(a + b)^2 = 1a^2b^0 + 2a^1b^1 + 1a^0b^2$$

$$(a + b)^3 = 1a^3b^0 + 3a^2b^1 + 3a^1b^2 + 1a^0b^3$$

and generally,

$$(a + b)^n = \binom{n}{0}a^n b^0 + \binom{n}{1}a^{n-1}b^1 + \binom{n}{2}a^{n-2}b^2 + \cdots + \binom{n}{n}a^0b^n$$

which we usually write using summation notation:

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

This result can be derived as follows: when we multiply out $(a + b)^n$, we have n factors $(a + b)$. We pick either an a or a b from each factor and multiply them all together to get a product. If the number of times we pick b is i , we get a term $a^{n-i}b^i$. We repeat this process using all possible choices, and add up the resulting terms. The number of terms $a^{n-i}b^i$ is just the number of ways of picking i of the b terms from the n factors—this is just the same as the number of ways of putting i balls in n boxes, and so is $\binom{n}{i}$. Thus, adding all the terms up, we get

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

There is a nice formula for the binomial coefficients:

$$\binom{n}{r} = \frac{n!}{r!(n-r)!}$$

This can be seen as follows: Label the r balls from 1 to r . There are n ways of putting the first ball into the n boxes. After this ball has been placed, there are $n - 1$ boxes free, so there are $n - 1$ ways of placing ball 2 *for each placement of ball 1*. Similarly, for each way of placing balls 1 and 2, there are $n - 2$ ways of placing ball 3, and so on. Thus, there are

$$n(n-1)(n-2)\cdots(n-r+1) = \frac{n!}{(n-r)!}$$

ways of placing the numbered balls in the boxes. Any particular set of r boxes, however, can contain $r!$ permutations of the r balls, which are indistinguishable if we erase the numbering on the balls. Hence to get the number of ways of placing r indistinguishable balls in n boxes, we have to divide this number by $r!$, yielding

$$\frac{n!}{r!(n-r)!}$$

Another way of thinking about this is to consider the equivalent problem of placing r white balls and $n - r$ black balls in n boxes. If we number the white balls and the black balls, then we have n numbered balls, which can be placed in the n boxes in $n!$ different ways. Then as before, we have to divide by $r!$ (for the white balls) and $(n - r)!$ (for the black balls).

If we look at the fifth row of Pascal's triangle, we see that all the terms except the outermost ones are divisible by 5. This is true in general in the p^{th} row, provided that p is a prime: all the numbers except the outermost ones will be divisible by p . We can see this as follows: the number

$$\binom{p}{i} = \frac{p!}{i!(p-i)!}$$

is an integer². If $0 < i < p$, then both i and $(p - i)$ are strictly less than p , and so neither factor in the denominator can contribute a factor of p . But the numerator contains a factor of p , and so this factor remains uncanceled in the fraction, which shows that it is a factor of $\binom{p}{i}$, and so $p \mid \binom{p}{i}$, which is what we wanted to show.

As a consequence, we have the striking result

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

for any prime p . (All the other terms in the binomial expansion on the right-hand side drop out because p divides all of their coefficients.)

In fact, the result remains true for powers of p :

1.5 Theorem *If p is a prime and if n is a positive integer, then*

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

PROOF. Again this is easy to prove by induction. We have already seen that it

²That the right hand side is an integer is not at all obvious. But we have shown that it is an entry in Pascal's triangle, and we know that all such entries are integers.

is true for $n = 1$. If it is true for n , then we can compute:

$$\begin{aligned} (a+b)^{p^{n+1}} &\equiv ((a+b)^{p^n})^p \pmod{p} && \text{(by definition)} \\ &\equiv (a^{p^n} + b^{p^n})^p \pmod{p} && \text{(by the inductive hypothesis)} \\ &\equiv (a^{p^n})^p + (b^{p^n})^p \pmod{p} && \text{(by the result for } n = 1) \\ &\equiv a^{p^{n+1}} + b^{p^{n+1}} \pmod{p} && \text{(by definition)} \end{aligned}$$

and we are done. □

1.3 Groups

1.3.1 Definition and examples

A *group* consists of the following:

- a set G (the elements, or members, of the group),
- a distinguished element (often denoted by e) of G . This element is called the *identity* element (or simply the identity) of G ,
- a binary operation (usually denoted by addition or multiplication) from $G \times G$ to G such that for all elements a, b and c of G ,
 1. $a(bc) = (ab)c$ (multiplication is associative)
 2. $ea = a = ae$ (e is a multiplicative identity)
- and such that to each element a of G there is a unique two-sided multiplicative inverse. That is, for each a there is a unique element b such that $ab = ba = e$. The multiplicative inverse of a is conventionally written a^{-1} .

If the group operation is written as addition instead of multiplication, the inverse of a is denoted by $-a$, and the identity element is often just denoted by 0 . Writing the group operation as multiplication, we can write repeated multiplication as exponentiation in the usual fashion:

$$g^n \stackrel{\text{def}}{=} g^{n-1} \cdot g$$

and we can also define

$$g^{-n} \stackrel{\text{def}}{=} (g^n)^{-1}$$

It is easy to see that this is also equal to $(g^{-1})^n$, and that the usual formulas remain true:

$$g^{n+m} = g^n \cdot g^m$$

$$(g^n)^m = g^{nm}$$

for all integers n and m .

Note that a set G could have two different group operations defined on it—in that case, we really have two different groups. (We will see an example of this below.) In other words, the group consists not only of the elements of G but also the group operation. Nevertheless, we will often simply refer to “the group G ”, with the understanding that the group operation is also implicitly being referred to. This sort of usage is called in mathematics an “abuse of notation”, and is a good thing as long as it does not lead to confusion. In fact, without abuse of notation, most mathematics would become completely unreadable.

Nevertheless, we may occasionally (but very rarely) want to be quite specific. In such a case, we may write something like “let $\langle G, \circ, e \rangle$ be a group”. The meaning of this is that G is the set of elements of the group, \circ is the group operation (so the group operation applied to the elements a and b of the group is $a \circ b$), and e is the identity element of the group.

Here are some examples of groups:

The real numbers \mathbf{R} , under addition The identity element is 0. The (additive) inverse of any element a is $-a$.

The non-zero real numbers \mathbf{R}^* , under multiplication The identity element is 1. The (multiplicative) inverse of any element a is $a^{-1} = 1/a$.

The positive real numbers \mathbf{R}^+ , under multiplication This is essentially the same as the previous example. Since the positive real numbers are themselves a group under multiplication, with the same identity element 1, they form a *subgroup* of the larger group \mathbf{R}^* .

The 2-dimensional plane \mathbf{R}^2 The operation is vector addition, and the identity element is the zero vector $(0, 0)$.

The 3-dimensional space \mathbf{R}^3 Similar to the previous example.

n -dimensional Euclidean space \mathbf{R}^n . Similar to the previous example. The elements of this space are n -tuples (x_1, x_2, \dots, x_n) . The group operation is (coordinate-wise) addition. The identity element is the zero vector $(0, 0, \dots, 0)$.

The integers \mathbf{Z} under addition The identity element is 0. (Note that the non-zero integers \mathbf{Z}^* do *not* form a group under multiplication, because no numbers except 1 and -1 have multiplicative inverses in \mathbf{Z}^* .)

The integers \mathbf{Z} mod m , under addition m can be any positive integer. This is a group for the following reasons:

As we noted above in Theorem 1.1, if $a_1 \equiv b_1 \pmod{m}$ and $a_2 \equiv b_2 \pmod{m}$, then $a_1 + a_2 \equiv b_1 + b_2 \pmod{m}$. This shows that it makes sense to talk about adding congruence classes (mod m); to add two congruence classes, just pick a number in each one and add them. The result will be in another congruence class which is by definition the sum of the first two congruence classes. The property of congruences noted above shows that it doesn't matter which two numbers you pick to add—you will always get the same final congruence class.

The congruence class containing 0 is then clearly the identity element, and the (additive) inverse of the congruence class containing a number a is the class containing the number $-a$. Thus, we really do have a group.

This group, which has m elements, is denoted \mathbf{Z}_m . As mentioned before, we often refer to the elements of this group as $\{0, 1, 2, \dots, m-1\}$, since any integer is congruent to exactly one of these integers (mod m).

The non-zero elements of \mathbf{Z}_p , under multiplication, for p prime If p is a prime, and a is not congruent to 0 (mod p), then p does not divide a . If also b is not congruent to 0 (mod p), then p does not divide the product ab . Thus the non-zero elements of \mathbf{Z}_p are closed under multiplication. The elements of this set are often denoted $\{1, 2, \dots, p-1\}$, as above, and the set itself is denoted \mathbf{Z}_p^* . Note that \mathbf{Z}_p^* has $p-1$ elements.

We have to show that \mathbf{Z}_p^* is a group. Just as we used Theorem 1.1 to show that addition of congruence classes is well-defined, we can use Theorem 1.2 to show—in exactly the same way—that multiplication of congruence classes is well-defined.

Clearly 1 is the multiplicative identity. So we have to show that every element has a multiplicative inverse.

To do this, we reason as follows: Let a be in \mathbf{Z}_p^* and consider the set obtained by multiplying a by each element of \mathbf{Z}_p^* . This set is

$$\{1 \cdot a, 2 \cdot a, 3 \cdot a, \dots, (p-1) \cdot a\}$$

Now these elements are all distinct, for to say that $ab = ac$ in \mathbf{Z}_p^* is just to say that $ab \equiv ac \pmod{p}$. But since $a \in \mathbf{Z}_p^*$, p is relatively prime to a , so this means that $b \equiv c \pmod{p}$; i.e., $b = c$ in \mathbf{Z}_p^* (by Theorem 1.4).

Since \mathbf{Z}_p^* is a finite set, it follows that the set of multiples of a is just a permutation of \mathbf{Z}_p^* . Therefore, one of the elements of this set must be 1.

That is, there must be an element $b \in \mathbf{Z}_p^*$ such that $ab \equiv 1 \pmod{p}$. This b is the multiplicative inverse of a .

For a concrete example, let us take $p = 5$ and $a = 2$. The set

$$\{1 \cdot 2, 2 \cdot 2, 3 \cdot 2, 4 \cdot 2\}$$

is just the set

$$\{2, 4, 1, 3\}$$

which is a permutation of

$$\mathbf{Z}_5^* = \{1, 2, 3, 4\}$$

This reasoning is used over and over in this subject.

The rotations in the plane around $(0, 0)$. The identity element is the rotation of 0 degrees, i.e., the “rotation” that does nothing. The inverse of the rotation through an angle θ is the rotation through the angle $-\theta$.

The rotations in \mathbf{R}^3 around lines passing through $(0, 0, 0)$ This is similar to the previous example.

All the groups listed above *except for the last one* have the property that the group operation is commutative. The last example is different, because rotations about different axes do not in general commute. Commutative groups are also called *abelian* groups in honor of Niels Henrik Abel, a Norwegian mathematician in the early 1800’s.

Although many important groups are non-abelian, we will only need to consider abelian groups in this survey.

1.3.2 Subgroups

A *subgroup* H of a group G is a subset of G that includes the identity of G and which is itself a group under the group operation of G . So for instance (as we pointed out already above), the multiplicative group \mathbf{R}^+ is a subgroup of the multiplicative group \mathbf{R}^* . Also, any line through the origin is an (additive) subgroup of \mathbf{R}^2 . Note, however, that a line that does not pass through the origin is *not* a subgroup of \mathbf{R}^2 : it does not contain the identity of \mathbf{R}^2 , and the vector sum of two elements of the line will not be on the line, so it is not even closed under the group operation.

Let us continue with this example. Let H be a line through the origin, thought of as a subgroup of \mathbf{R}^2 . If a vector a is not in the line H then $a + H$ (remember we are using additive notation here) is a line that is parallel to H . (It is called

a *translation* of H ; the term “translation” refers geometrically to a rigid motion not involving a rotation.) $a + H$ does not pass through the origin, so it is not a subgroup of \mathbf{R}^2 . And since $a + H$ is parallel to H , it is disjoint from H .

Now we could continue this process: if b is a vector that is not in H and also not in $a + H$ then $b + H$ is another line that is parallel to H and disjoint from both H and $a + H$, and so on. In fact, the whole group \mathbf{R}^2 can be thought of as being “tiled” by disjoint copies (or translates) of H , each copy being of the form $a + H$. (See Figure 1.1)

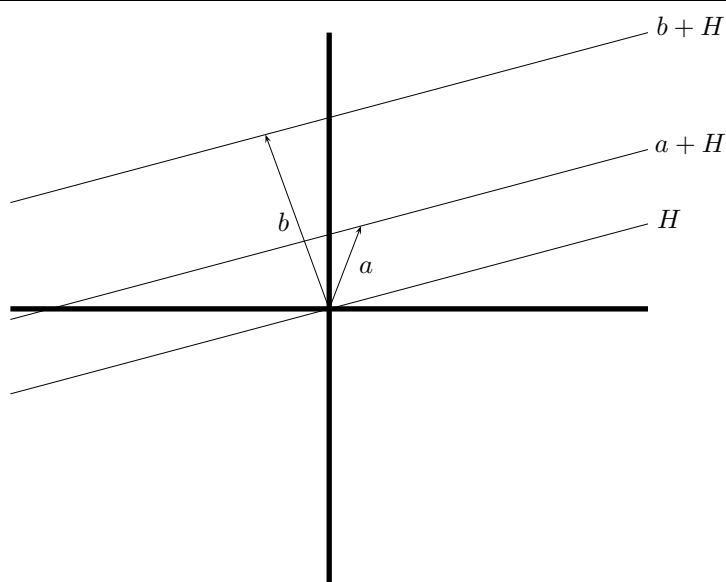


Figure 1.1: Tiling \mathbf{R}^2 with translates of a subgroup H

This construction is valid in any group: Let G be a group, with the group operation written as multiplication. If H is a subgroup of G , and if a is any element of G , then the set aH is defined to be the set of multiples

$$aH \stackrel{\text{def}}{=} \{ah : h \in H\}$$

You can think of the set aH as “ H translated by a ”.

If a is not in H , then aH and H have no elements in common. This is because if ah_1 were an element of aH that was also an element h_2 of H , then from the equation

$$ah_1 = h_2$$

we would get (by multiplying on the right by h_1^{-1})

$$ah_1h_1^{-1} = h_2h_1^{-1}$$

But the left hand side is just a , and the right hand side is an element of H , so this says that a is an element of H , which we assumed was not the case. Thus aH and H must be disjoint.

This shows again that aH is *not* a subgroup of G —for one thing, it doesn't contain the identity.

We can continue as before: if b is not in H and not in aH then bH is disjoint from H by what we have already shown. It is also disjoint from aH , for if $bh_1 = ah_2$ then just as before, we get $b = ah_2h_1^{-1}$, which shows that b is in aH , a contradiction. Thus, continuing the process until we can't go any further, we see that the whole group G is a disjoint union of sets of the form aH .

If G is a finite group, then H is also finite. Further, the number of elements in aH (for any a) is the same as the number of elements in H . For $ah_1 = ah_2$ if and only if $h_1 = h_2$ (by multiplying or dividing on the left by a). Thus, the number of elements of G is just the product of the following two numbers:

- the number of elements of H
- the number of distinct sets aH

and in particular, the number of elements of H divides the number of elements of G .

The number of elements of a group is called the *order* of that group; similarly for a subgroup. We have thus shown the following:

1.6 Theorem *The order of any subgroup of a finite group divides the order of the group.*

This theorem is due to Lagrange.

1.3.3 Products of groups

If $\langle G, \circ, e \rangle$ and $\langle H, \bullet, f \rangle$ are groups, then we can form a new group $G \times H$ simply by making each original group a “coordinate” in the new group and having the group operations act separately on each coordinate. To be precise, the new group consists of all the ordered pairs (g, h) where g is an element of G and h is an element of H . The group operation is defined by

$$(g_1, h_1) \times (g_2, h_2) \stackrel{\text{def}}{=} (g_1 \circ g_2, h_1 \bullet h_2)$$

With this definition, it is clear that the element (e, f) is the identity for the new group and that the inverse of (g, h) is (g^{-1}, h^{-1}) , where g^{-1} is the inverse (in G) of g , and similarly for h^{-1} .

You may be concerned that, while it is clear that (e, f) is an identity, we have not shown that the identity is unique. Well, this is an easy exercise. But actually, it is not necessary to prove this, because in a complete treatment of group theory, one of the first results proved is that if an identity exists at all, it is automatically unique. (The proof is quite easy.) So we don't have to concern ourselves with this point.

Products of groups are quite common. For instance, $\mathbf{R} \times \mathbf{R}$ is just the group \mathbf{R}^2 . And $\mathbf{R} \times \mathbf{R} \times \mathbf{R}$ (it's obvious what this means—the group of ordered triples. . .) is just the group \mathbf{R}^3 . And so on.

Here are some other important examples:

$\mathbf{Z}_2 \times \mathbf{Z}_3$ This group has 6 elements, which we can write as

$$\begin{array}{ccc} (0, 0) & (0, 1) & (0, 2) \\ (1, 0) & (1, 1) & (1, 2) \end{array}$$

where $(0, 0)$ is the identity, and addition is (mod 2) in the first coordinate and (mod 3) in the second.

The group \mathbf{Z}_6 also has 6 elements, and is *isomorphic* to $\mathbf{Z}_2 \times \mathbf{Z}_3$. This means that each group can be thought of as a simple renaming of the other. Precisely, two groups $\langle G, \circ, e \rangle$ and $\langle H, \bullet, f \rangle$ are isomorphic if and only if there is a function $\phi : G \rightarrow H$ (the “renaming function”) such that

- ϕ is 1-1 and onto (i.e., ϕ sets up a 1-1 correspondence between the elements of G and the elements of H),
- $\phi(e) = f$, and
- For all elements a and b of G , $\phi(a \circ b) = \phi(a) \bullet \phi(b)$

We can set up the correspondence between $\mathbf{Z}_2 \times \mathbf{Z}_3$ and \mathbf{Z}_6 (which is represented by the set $\{0, 1, 2, 3, 4, 5\}$) as follows:

$$\begin{array}{ccc} (0, 0) & \longleftrightarrow & 0 \\ (0, 1) & \longleftrightarrow & 4 \\ (0, 2) & \longleftrightarrow & 2 \\ (1, 0) & \longleftrightarrow & 3 \\ (1, 1) & \longleftrightarrow & 1 \\ (1, 2) & \longleftrightarrow & 5 \end{array}$$

In terms of the renaming function ϕ , this is just the same as defining³

$$\begin{aligned}\phi(0,0) &= 0 \\ \phi(0,1) &= 4 \\ \phi(0,2) &= 2 \\ \phi(1,0) &= 3 \\ \phi(1,1) &= 1 \\ \phi(1,2) &= 5\end{aligned}$$

It is a straightforward exercise to show that this correspondence preserves the group operation. (That's very important—just a 1-1 correspondence is completely insignificant unless it also preserves the group operation.) We'll explain below an easy way to come up with this correspondence.

$\mathbf{Z}_2 \times \mathbf{Z}_2$ This is a group with four elements. One might think at first that since $\mathbf{Z}_2 \times \mathbf{Z}_3$ is isomorphic to \mathbf{Z}_6 , that also $\mathbf{Z}_2 \times \mathbf{Z}_2$ is isomorphic to \mathbf{Z}_4 . Certainly they have the same number of elements, so they can be put into 1-1 correspondence. But there is no way to do this so that the group operation is preserved. We will see why this is true in the next subsection.

1.3.4 The order of an element of a group

Let G be a finite group, with the group operation represented as multiplication. If g is any element of G , consider the set of powers $\{g^i : i = 1, 2, \dots\}$ of g . Since g is a finite group, eventually there will be two powers that are equal—say they are g^c and g^d . We may assume that $c > d$. Then since $g^c = g^d$, we have $e = g^c(g^d)^{-1} = g^c g^{-d} = g^{c-d}$. Thus, there is a positive integer n (in this case $n = c - d$) such that $g^n = e$.

The *order* of an element g of G is defined to be the least positive integer n such that $g^n = e$. (Note that we have previously defined the order of a group. These two uses of the word “order” are not precisely the same, but they are related, as we'll see below.)

Here are some simple things we can say about the order of an element. Say n is the order of g :

- g raised to any multiple of n is e . For, $g^{rn} = (g^n)^r = e^r = e$.
- If $g^m = e$ then m is a multiple of n . For if not, we can divide m by n in the usual way, to get $m = qn + r$ where the remainder r is not 0; we have $0 < r < n$. Then $e = g^m = g^{qn+r} = (g^n)^q g^r = g^r$, and this is a contradiction, because $0 < r < n$.

³Note that actually we should write $\phi((0,0))$ instead of $\phi(0,0)$. But since there is no chance of confusion here, we use the shorter form. This is another example of abuse of notation.

So the powers of g that yield the identity are precisely the multiples of the order of g .

Already we can now see that $\mathbf{Z}_2 \times \mathbf{Z}_2$ cannot be isomorphic to \mathbf{Z}_4 . For every element of $\mathbf{Z}_2 \times \mathbf{Z}_2$ has order 1 or 2, while two elements of \mathbf{Z}_4 (1 and 3) have order 4. So these two groups, even though they have the same number of elements, cannot be thought of as just renamings of each other—the group operations act in very different ways on the two groups.

To repeat: if two groups are isomorphic, they must have the same number of elements. But just because two groups have the same number of elements does not guarantee that they are isomorphic. The group operation must also be preserved by the renaming map.

A finite group is *cyclic* if there is an element whose order is the order of the group. Another way of saying this is that a group is cyclic if there is an element g such that every element of the group is a power of g . In such a case, we say that such an element g *generates* the group. What we pointed out above is that \mathbf{Z}_4 is cyclic, but $\mathbf{Z}_2 \times \mathbf{Z}_2$ is not, so they cannot be isomorphic.

The way we arrived at the correspondence between $\mathbf{Z}_2 \times \mathbf{Z}_3$ and \mathbf{Z}_6 was to note that both groups are cyclic. We picked a generator for each, and made them correspond. (The generators we picked were 1 for \mathbf{Z}_6 and $(1, 1)$ for $\mathbf{Z}_2 \times \mathbf{Z}_3$.) Then corresponding powers⁴ of these elements must also correspond; that determines the renaming map. Actually, \mathbf{Z}_6 has one other generating element: -1 also has order 6. We could have picked either 1 or -1 to use in constructing the isomorphism.

Clearly, all the additive groups \mathbf{Z}_m are cyclic, and they are all generated by the element 1. It is a remarkable fact that the multiplicative groups \mathbf{Z}_p^* are also cyclic. I don't think this is at all obvious. The proof, however, is stunningly short—we will see it later (Theorem 1.18, page 42).

If h is any element of a finite group G , let H denote the set of all powers of h . So H is all of G if and only if G is cyclic and h generates G . But even if this is not true, H itself is a cyclic subgroup of G generated by h . And clearly the order of h is the number of elements of the subgroup H . By Theorem 1.6 (page 17), then, we have:

1.7 Theorem *The order of any element of a finite group divides the order of the group.*

Let's see what this means in two simple cases:

$G = \mathbf{Z}_m$ If a is an element of \mathbf{Z}_m (i.e., a number (mod m)), then the order n of a divides m . To say that n is the order of a just says that $na \equiv 0 \pmod{m}$

⁴actually, corresponding multiples, since we are using the additive notation for the group operations

m), and as we have seen, the fact that $n|m$ means that also $ma \equiv 0 \pmod{m}$. This is really trivial—of course $m|ma$ for any number a .

$G = \mathbf{Z}_p^*$, for p prime The order of G is $p - 1$. If n is the (multiplicative) order of a number a in \mathbf{Z}_p^* , then $a^n \equiv 1 \pmod{p}$. The theorem tells us that n divides $p - 1$, and so also $a^{p-1} \equiv 1 \pmod{p}$.

In contrast to the result for \mathbf{Z}_m , the result for \mathbf{Z}_p^* is not at all trivial—I don't think it's obvious at all, even though the proof we just gave is today quite simple. It is significant enough to have a name:

1.8 Theorem (Fermat's "little theorem") *If p is a prime and p does not divide a , then*

$$a^{p-1} \equiv 1 \pmod{p}$$

1.3.5 Orders of elements in abelian groups

1.9 Lemma *If a and b are two elements in an abelian group whose orders are relatively prime, then the order of ab is the product of their orders.*

PROOF. Say the order of a is n and the order of b is m . Since a and b commute, we have

$$(ab)^{nm} = a^{nm}b^{nm} = e$$

so the order of ab is at most nm .

Now say the order of ab is r . We have

$$(ab)^{rn} = \begin{cases} ((ab)^r)^n & = e \\ (a^n)^r b^{rn} & = b^{rn} \end{cases}$$

so $b^{rn} = e$. Therefore the order of b , which is m , must divide rn . But n and m are relatively prime. Therefore m actually divides the order r of ab .

Similarly, n must also divide r . Hence $r \geq \text{lcm}(m, n) = mn$ □

Now one might think at first that one could extend this theorem so that in general (even if the orders of a and b were not relatively prime), the order of ab would be $\text{lcm}(n, m)$. However, this is clearly not true. For instance, just take $b = a^{-1}$. Then a and b have the same order, but the order of ab is 1. Nevertheless, we can show that there is *some* element of the group whose order is $\text{lcm}(n, m)$:

1.10 Lemma *If an element a of an abelian group G has order n , and if d is any divisor of n , then there is an element of G of order d .*

PROOF. $a^{n/d}$ is such an element. For $(a^{n/d})^d = a^n = e$, so the order of $a^{n/d}$ is no larger than d . And on the other hand, its order can't be less than d , because if it were, then the order of a would be less than n , which is a contradiction. \square

1.11 Theorem *If a and b are two elements in an abelian group G whose orders are n and m respectively, then there is an element c of G whose order is $\text{lcm}(n, m)$.*

PROOF. Say n and m have the prime factorizations

$$\begin{aligned} n &= \prod_{p_i \text{ prime}} p_i^{\nu_i} \\ m &= \prod_{p_i \text{ prime}} p_i^{\mu_i} \end{aligned}$$

Then as we have already mentioned,

$$\text{lcm}(n, m) = \prod_{p_i \text{ prime}} p_i^{\max(\nu_i, \mu_i)}$$

For each prime p_i , we know by the proof of the preceding lemma that

$$a^{n/p_i^{\nu_i}} \text{ has order } p_i^{\nu_i}$$

and similarly

$$b^{m/p_i^{\mu_i}} \text{ has order } p_i^{\mu_i}$$

Set

$$c_i = \begin{cases} a^{n/p_i^{\nu_i}} & \text{if } \nu_i \geq \mu_i \\ b^{m/p_i^{\mu_i}} & \text{otherwise} \end{cases}$$

Then c_i has order $p_i^{\max(\nu_i, \mu_i)}$. The orders of all the elements c_i are relatively prime, so their product $c = \prod c_i$ has order $\prod p_i^{\max(\nu_i, \mu_i)} = \text{lcm}(n, m)$. \square

1.12 Corollary *If G is an abelian group, and if n is the greatest order of any element of G , then the order of each element of G divides n .*

PROOF. If m is the order of an element of G , and if m does not divide n , then there is an element of order $\text{lcm}(m, n)$, which must be strictly greater than n . But we assumed that n was the largest order of any element of G , so this is a contradiction. \square

1.4 Rings

Well, groups are quite nice, but the mathematical objects we know and love best really have *two* operations, not one. Rings and fields are abstract structures having two operations. These operations are almost always denoted by addition and multiplication. You can think of a ring as an abstract version of the integers \mathbf{Z} , and a field as an abstract version of the real numbers \mathbf{R} .

A *ring* is a set R , together with two operations, denoted by addition and multiplication, such that

- Under addition, R is an abelian group. The identity is denoted by the symbol 0 . Thus, for all elements a and b of R ,

$$a + b = b + a$$

$$a + 0 = a = 0 + a$$

- Multiplication is associative, and there is a multiplicative identity element, denoted by the symbol 1 . Thus,

$$a(bc) = (ab)c$$

$$a1 = a = 1a$$

- Multiplication distributes over addition. That is, for all elements a , b , and c of R ,

$$a(b + c) = ab + ac$$

$$(b + c)a = ba + ca$$

We denote the element $1 + 1$ by 2 . Similarly, we denote $1 + 2$ by 3 , and so on for any positive integer n . The additive inverse of n is (as usual) denoted by $-n$. Thus, for any element a in the ring,

$$2a = (1 + 1)a = 1a + 1a = a + a$$

and in general, for any integers n and m ,

$$(n + m)a = na + ma$$

Certainly the ordinary integers \mathbf{Z} form a ring (with the ring elements 0 and 1 being the actual integers 0 and 1). In fact, the integers have a few more properties, which we will need:

- Multiplication is commutative: for elements a and b ,

$$ab = ba$$

A ring in which multiplication is commutative is called a *commutative ring*. (We have already defined a ring so that *addition* is automatically commutative.)

Finally,

- If a and b are two elements and $ab = 0$, then at least one of a and b is 0.

This is usually expressed by saying that the ring has no zero-divisors. It may not at first be obvious that a ring might have zero-divisors, because we are so used to thinking of the ordinary integers, where this can't happen. However, here is an example:

The set \mathbf{Z}_6 is a ring under addition and multiplication. This is because just as addition of congruence classes is well-defined by Theorem 1.1, multiplication of congruence classes is well-defined by Theorem 1.2. However, we see that in this ring, $2 \cdot 3$ is 0. (That is, $2 \cdot 3 \equiv 0 \pmod{6}$.) So this ring has zero-divisors.

On the other hand, if p is a prime, then the ring \mathbf{Z}_p does not have zero-divisors. This is because, as we have seen already before, if ab is (congruent to) $0 \pmod{p}$, then $p|ab$, so either $p|a$ or $p|b$; i.e., either a or b is already (congruent to) $0 \pmod{p}$.

A commutative ring with no zero-divisors is called an *integral domain*, presumably because it is very similar to the ordinary integers. All the rings we will need to consider in this survey are integral domains.

Note that although a ring is a group under addition, its non-zero elements do not necessarily constitute a group under multiplication—although it has a multiplicative identity, it does not have to have multiplicative inverses. The ring \mathbf{Z} is an example.

For now, the two main examples of rings to keep in mind are \mathbf{Z} and \mathbf{Z}_p , where p is any prime. Soon we shall see some more rings, which will be very important to us.

1.5 Fields

A *field* is a commutative ring in which each non-zero element has a multiplicative inverse.

If R is a ring, we denote the set of non-zero elements of R by R^* . Thus, a commutative ring R is a field if R^* is a multiplicative group. (And in general, the set of non-zero elements of a field is called the *multiplicative group* of the field.)

Thus when R is a field, the set R^* is in particular closed under multiplication, and so R has no zero-divisors. That is, each field is automatically also an integral domain.

The ring \mathbf{Z} is not a field, as we have already noted. On the other hand, the ring \mathbf{Z}_p is a field, for any prime p —we saw previously that \mathbf{Z}_p^* was a group.

The most intuitive fields for most of us are the fields \mathbf{R} (the real numbers), \mathbf{C} (the complex numbers), and \mathbf{Q} (the rational numbers). However, even though these fields are the most important fields in mathematics, they are relatively unimportant for our purposes in this survey. The most important field for us is \mathbf{Z}_p . Note that \mathbf{Z}_p has a finite number of elements, and is thus the first example we have seen of a *finite field*. We will see that there are many more.

Just as a general group is often denoted by the letter G , a general field is most often denoted by K or k . This comes from the German word for field, *Körper*. (In French, similarly, the word is *corps*.)

We mentioned above that we can consider each integer n as being an element of any ring; this holds in particular then also for fields. It may be, however, that there are some positive integers that are equal to 0 in the ring. For instance, in \mathbf{Z}_p , any multiple of p is 0. In such a case, we call the smallest positive integer that is 0 in the ring the *characteristic* of the ring. So for instance, the characteristic of the ring \mathbf{Z}_m is m . It is easy to see that any integer then is 0 in the ring if and only if that integer is a multiple of the characteristic. (In fact, the characteristic is just the order of the element 1 in the additive group of the field, and so $n \cdot 1 = n = 0$ if and only if the order of 1 divides n .)

If no positive integer equals 0 in the ring, you might think we would say that the ring has characteristic infinity. However, algebraists say in such a case that the ring has characteristic 0. So for instance, the rings \mathbf{Z} , \mathbf{R} , and \mathbf{Q} all have characteristic 0. However, all finite fields must have non-zero characteristic, since, being finite, they can't contain distinct copies of *all* the integers.

The characteristic of any field must be prime. This is because otherwise, there would be two non-zero numbers whose product was the characteristic; i.e., the field would have zero-divisors, and we know this is impossible.

In a field K of characteristic p , the elements $\{0, 1, 2, \dots, p-1\}$ themselves form a field that is obviously isomorphic to \mathbf{Z}_p . This subfield of K is called the *prime field* of K .

The characteristic of a finite field divides the number of elements in the field. This is simply because the prime field is a subgroup of the additive group of the field, and the order of the subgroup must divide the order of the group. This result is actually much weaker than the whole truth, which is that the number of elements in a finite field is a *power* of the characteristic. We'll see that below in Theorem 1.17 on page 42.

1.6 Vector spaces

1.6.1 Vector spaces and dimensionality

The simplest vector spaces are \mathbf{R}^2 and \mathbf{R}^3 . When we call them vector spaces, we are referring to the following properties, which we will illustrate using the vector space \mathbf{R}^2 :

- Their elements are called *vectors*. A typical vector is $\vec{v} = (2, -3)$.
- Vectors can be added. Addition is commutative, and the set of vectors forms a group under addition. That is, there is an additive identity (the zero vector $(0, 0)$), and every vector \vec{v} has an additive inverse $-\vec{v}$. For instance, $-(2, -3)$ is the vector $(-2, 3)$.
- The real numbers are called *scalars*. Vectors can be multiplied by scalars to form other vectors. For instance, $5(2, -3) = (10, -15)$. This scalar multiplication “acts like you would expect”. For instance,

$$-(2, -3) = (-1)(2, -3)$$

$$(7 + \pi)(2, -3) = 7(2, -3) + \pi(2, -3)$$

$$7((2, -3) + (0, 1/3)) = 7(2, -3) + 7(0, 1/3)$$

and so on.

We call \mathbf{R}^2 a vector space over \mathbf{R} , because the scalars are elements of \mathbf{R} . \mathbf{R}^2 can be represented as the set of all ordered pairs of real numbers.

Similarly, \mathbf{R}^3 is also a vector space over \mathbf{R} , and can be represented as the set of all ordered triples of real numbers.

In general, we can make the following definition:

A *vector space over \mathbf{R}* is a set V such that

- V is an additive abelian group. The additive identity is denoted by $\vec{0}$ and is called the zero vector.
- There is an operation called *scalar multiplication* that takes an element of \mathbf{R} and an element of V and produces another element of V . Scalar multiplication is written multiplicatively, and has the following properties: for all vectors \vec{v} and \vec{w} and for all elements a and b of \mathbf{R} ,

1. $a(b\vec{v}) = (ab)\vec{v}$

2. $(a + b)\vec{v} = a\vec{v} + b\vec{v}$
3. $a(\vec{v} + \vec{w}) = a\vec{v} + a\vec{w}$
4. $1\vec{v} = \vec{v}$
5. $0\vec{v} = \vec{0}$ (This actually follows from item 2.)
6. $(-1)\vec{v} = -\vec{v}$ (This follows from items 2, 5, and 4.)

We can go one step farther, and for the purposes of this survey, this is important: There is no reason why the scalars have to be elements of \mathbf{R} . All that is really needed is that there be a field K which can be used as the scalar field. A vector space over a field K is defined exactly as above, simply substituting K for \mathbf{R} . Vector spaces over the complex numbers \mathbf{C} are important in quantum mechanics and in many parts of mathematics. For us, a useful example is a vector space over \mathbf{Z}_p : pick a prime p and consider all ordered pairs of elements of \mathbf{Z}_p . Scalar multiplication is performed coordinate-wise; that is,

$$a(x, y) = (ax, ay)$$

It is easy to see that this is actually a vector space over \mathbf{Z}_p . For instance—say $p = 3$ —we have

$$(2, 1) + (1, 1) = (0, 2)$$

$$2(2, 1) = (1, 2)$$

(Remember that all computations are carried out in \mathbf{Z}_3 , that is, (mod 3)).

This vector space, which is typically denoted \mathbf{Z}_p^2 is a *2-dimensional* vector space over \mathbf{Z}_p , because it can be represented as ordered *pairs* of elements of \mathbf{Z}_p . Similarly, \mathbf{Z}_p^n , the set of ordered *n-tuples* of elements of \mathbf{Z}_p , is an *n-dimensional* vector space over \mathbf{Z}_p .

In general, if K is a field, K^n denotes the vector space of *n-tuples* of elements of K . If K is a finite field—say K has c elements—then K^n will have c^n elements, since there will be c^n different ordered pairs with coordinates in K . For instance, \mathbf{Z}_p^n has p^n elements.

Sometimes a set can be considered to be a vector space over 2 different fields. In this case, the dimension of the space depends on the field over which it is a vector space. For instance, the vector space \mathbf{C}^2 is the 2-dimensional space of ordered pairs of complex numbers:

$$\mathbf{C}^2 = \{(z_1, z_2) : z_1, z_2 \in \mathbf{C}\}$$

Now each complex number z has a real and an imaginary part—we write $z = x + iy$, where x and y are each real numbers. Thus, we could write

$$\mathbf{C}^2 = \{(x_1 + iy_1, x_2 + iy_2) : x_1, y_1, x_2, y_2 \in \mathbf{R}\}$$

which is isomorphic to the set of all *real* 4-tuples:

$$\mathbf{R}^4 = \{(x_1, y_1, x_2, y_2) : x_1, y_1, x_2, y_2 \in \mathbf{R}\}$$

Thus, as a vector space over the field \mathbf{C} , \mathbf{C}^2 has dimension 2. But as a vector space over the field \mathbf{R} , \mathbf{C}^2 has dimension 4.

Note that these are really two different vector spaces, even though they consist of the same set of elements: when we consider \mathbf{C}^2 as a vector space over \mathbf{R} , we can only multiply each vector by real numbers; when we consider it as a vector space over \mathbf{C} , we can multiply each vector by any complex number.

As a matter of notation, if we have a vector space of n -tuples over a field K , we give names to the following special vectors:

$$\begin{aligned} e_1 &= \langle 1, 0, \dots, 0 \rangle = \text{“unit vector in the } x \text{ direction”} \\ e_2 &= \langle 0, 1, \dots, 0 \rangle = \text{“unit vector in the } y \text{ direction”} \\ &\vdots \\ e_n &= \langle 0, 0, \dots, 1 \rangle = \text{“unit vector in the } n^{\text{th}} \text{ direction”} \end{aligned}$$

All the vector spaces we have been looking at have been sets of ordered n -tuples of elements of a field K . Such a space is said to have dimension n over K .

A set of vectors $\{v_1, v_2, \dots, v_r\}$ is said to *generate* or *span* the vector space V if every element v of V can be represented as a linear combination of elements v_i with coefficients in K :

$$v = k_1v_1 + k_2v_2 + \dots + k_nv_n$$

(Of course some or all of the coefficients k_i might be 0.) For instance, the vectors $\{e_1, e_2, \dots, e_n\}$ span the vector space K^n of n -tuples of elements of the field K .

One of the main theorems proved early on in courses in linear algebra is this:

1.13 Theorem *If V is a vector space over a field K for which there is a set of finitely many vectors that generate all of V over K , then V is actually finite-dimensional—that is, V can be represented as the set of all n -tuples of elements of K . Further, the dimension n is unique—one could not have one representation of V as the set of 3-tuples over K and another as the set of 4-tuples over K .*

We won't prove this theorem here, but we do use it for one important observation. Remember that we showed above on page 25 that the characteristic of a finite field divides the number of elements of the field. That's really only part of the story; the truth of the matter is quite a bit stronger:

1.14 Theorem *If K is a finite field of characteristic p , then the number of elements of K is a power of p .*

PROOF. Just by referring to the definition above of a vector space, we see that any field K is a vector space over its prime field. Further, any finite field K , when considered as a vector space over its prime field, must have a finite basis, since there are only finitely many elements of K to begin with. Therefore K is a finite-dimensional vector space over its prime field. Say its dimension is n . Then K can be represented as the set of n -tuples of elements of the prime field. Since there are p^n such n -tuples, it follows that a finite field of characteristic p has p^n elements. \square

Theorem 1.13 is actually usually expressed somewhat differently: We say that a set of vectors $\{v_1, v_2, \dots, v_r\}$ that spans a vector space V is a *basis* if the set is minimal—that is, if there is no subset of this set that also spans V . Then the theorem states that each basis of a vector space V has the same number of elements. This number is then called the dimension of V .

Continuing a little with this, a set of vectors $\{v_j\}$ is *linearly independent* if no vector in this set can be expressed as a linear combination of the remaining vectors. If $\{v_j\}$ is a spanning set for V and it is not linearly independent, then some vector (say v_1) can be expressed in terms of the rest of the vectors; we would have

$$v_1 = a_2v_2 + a_3v_3 + \cdots + a_rv_r$$

But this means that any vector v in V , which by assumption has a representation in terms of the vectors $\{v_1, v_2, \dots, v_r\}$:

$$v = b_1v_1 + b_2v_2 + b_3v_3 + \cdots + b_rv_r$$

actually has a representation in terms of the smaller set $\{v_2, \dots, v_r\}$ —we can just rewrite v_1 in terms of the remaining vectors as above and collect the resulting terms. This shows that if a spanning set is not linearly independent, it is not a basis of V .

Conversely, if a spanning set $\{v_1, v_2, \dots, v_r\}$ is not a basis of V , then it has a subset which still spans V . If this subset does not include v_1 , say, then this just means that v_1 (as an element of V) must be a linear combination of the remaining vectors in this set, and so the set is not linearly independent.

So a set of vectors $\{v_j\}$ is a basis if and only if it spans V and is linearly independent.

In case V has dimension 2, any basis will have 2 elements. In this case, things are particularly simple: two vectors are linearly independent if and only if neither is a multiple of the other. Figure 1.2 shows an example of two bases for \mathbf{R}^2 , considered as a vector space over \mathbf{R} . The first basis is the usual one:

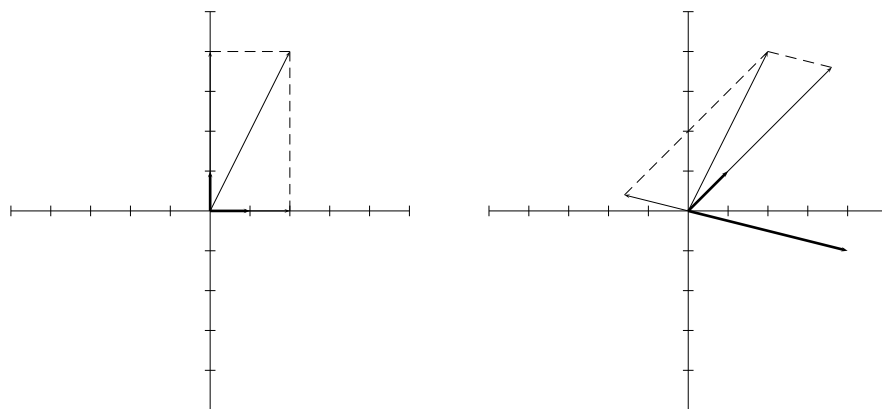
$$\begin{aligned} e_1 &= (1, 0) \\ e_2 &= (0, 1) \end{aligned}$$

The second basis consists of the two vectors

$$f_1 = (4, -1)$$

$$f_2 = (1, 1)$$

Clearly in each case, neither vector is a multiple of the other. The figure shows how the vector $(2, 4)$ can be expressed as a linear combination of the basis vectors in each case.



$$(2, 4) = 2e_1 + 4e_2$$

$$(2, 4) = -\frac{2}{5}f_1 + \frac{18}{5}f_2$$

Figure 1.2: Two different bases for \mathbf{R}^2 . In each case the basis vectors are the thick ones.

1.6.2 Linear functions

If V is a vector space over a field K , a function $L : V \rightarrow V$ (“ L maps V to V ”) is said to be a *linear function* (or a *linear transformation*) if and only if the following two criteria both hold:

1. $L(u + v) = L(u) + L(v)$ for all vectors u and v in V .
2. $L(au) = aL(u)$ for all vectors u in V and scalars a in K .

We could combine these conditions into the single condition

1. $L(au + bv) = aL(u) + bL(v)$ for all vectors u and v and scalars a and b .

It follows immediately (take $u = v$, $a = 1$, and $b = -1$) that for any linear function L , $L(0)$ must be the zero vector.

If V is an n -dimensional vector space over K , represented as the set of ordered n -tuples over K , then any $n \times n$ matrix with elements in K corresponds to a linear function from V to V , as follows: say the matrix is

$$T = \begin{pmatrix} t_{11} & t_{12} & \cdots & t_{1n} \\ t_{21} & t_{22} & \cdots & t_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ t_{n1} & t_{n2} & \cdots & t_{nn} \end{pmatrix}$$

The matrix T applied to the vector $v = \langle v_1, v_2, \dots, v_n \rangle$ is just the vector whose i^{th} coordinate is

$$\sum_{j=1}^n t_{ij}v_j$$

That is, the i^{th} coordinate of the vector $T(v)$ is the “dot product” of the i^{th} row of the matrix T with the vector v .

It is easy to see that this function is linear. It is also true that *any* linear function from V to V has this form—it can be represented by a matrix in exactly this way. We don’t actually need this result here, although it is easy to prove.

What is important for us is that composition of linear functions corresponds to multiplication of their associated matrices. That is, if S and T are linear functions given by the matrices $\{s_{ij}\}$ and $\{t_{ij}\}$ then the linear function ST (“first apply T , then apply S to the result”), which is defined so that

$$(ST)(v) = S(T(v))$$

is given by a matrix $A = \{a_{ij}\}$, where

$$a_{ij} = \sum_{k=1}^n s_{ik}t_{kj}$$

If T is a matrix, we define $T^2 = TT$, $T^3 = TTT$, and so on. In this way, we can define T^n for any positive integer n .

The identity matrix I is defined to be the matrix that is 1 on the diagonal and zero elsewhere. That is, if I is the matrix a_{ij} , then

$$a_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}$$

I is then the identity function, when considered as a linear transformation: $I(v) = v$ for every v in V . Since $IT = TI = T$ for any matrix T , it makes sense to write $T^0 = I$ for any T .

We should also note that if T is the 0 transformation (i.e., if $Tv = 0$ for all vectors v)⁵ then the matrix corresponding to T must be the 0 matrix. This is because Te_j is just the vector whose components are the j^{th} column of this matrix⁶. So if all Te_j are 0 then every element in every column of the matrix for T must be 0.

In a similar way, we can show that a linear transformation is determined by its values on a set of vectors that spans V :

Say $\{w_i\}$ is a set of vectors⁷ in V that spans V . That is, every vector v in V can be written as a linear combination of vectors $\{w_i\}$ with coefficients in K . Suppose that we have a linear transformation T , and we know that $T(w_i) = 0$ for each vector w_i . Then T is identically 0. This is simply because if v is any vector, we can write

$$v = a_1w_1 + a_2w_2 + \cdots + a_nw_n$$

with the coefficients a_i in K . Then because T is linear, we get

$$T(v) = a_1T(w_1) + a_2T(w_2) + \cdots + a_nT(w_n)$$

and each term on the right is 0, since all the $T(w_i)$ are 0. So $T(v) = 0$ for all vectors v , and so T is identically 0.

The way this result is used is that if we have two matrices A and B , and we know that $Aw_i = Bw_i$ for all i , then we must have $A = B$. For we can just set $T = A - B$ and apply the previous reasoning.

1.6.3 Dot products and adjoint transformations

We mentioned dot products above. We can define a dot product (or “scalar product”) of two vectors in the usual fashion. Usually mathematicians do not write it with a dot, however, but with parentheses, like this:

$$(u, v) = \sum_{i=1}^n u_i v_i$$

⁵As a slight abuse of notation we often write Tv instead of $T(v)$. The meaning is always clear from the context.

⁶Remember that $e_1 = \langle 1, 0, \dots, 0 \rangle$, $e_2 = \langle 0, 1, \dots, 0 \rangle$, and so on.

⁷This is a little bit tricky notationally: I really wanted to write $\{v_i\}$, but I have used v_i above to denote a (scalar) component of the vector v . So here I am using w_i , with the understanding that w_i is a vector in V , not a scalar component in K .

Of course this looks like an ordered pair, but it is not. You have to get the meaning by context, but this is never a problem in practice. Note that here u and v are vectors in V , but (u, v) is a scalar in K .

If T is a linear transformation from V to V , as above, and we represent T as a matrix $\{t_{ij}\}$, let us denote by T^* the transformation whose matrix is the transpose of T . That is,

$$t_{ij}^* = t_{ji}$$

Then given any two vectors u and v in V , we have

$$\begin{aligned} (Tu, v) &= \sum_{i=1}^n (Tu)_i v_i \\ &= \sum_{i=1}^n \sum_{j=1}^n t_{ij} u_j v_i \\ &= \sum_{j=1}^n u_j \sum_{i=1}^n t_{ij} v_i \\ &= \sum_{j=1}^n u_j \sum_{i=1}^n t_{ji}^* v_i \\ &= (u, T^*v) \end{aligned}$$

T^* is called the *adjoint* of T . Note that since

$$(STu, v) = (Tu, S^*v) = (u, T^*S^*v)$$

we have

$$(ST)^* = T^*S^*$$

(note the order reverses), and consequently

$$(T^n)^* = (T^*)^n$$

1.6.4 An application of vector spaces

The familiar Fibonacci sequence

$$1, 1, 2, 3, 5, 8, 13, 21, 34, 55, 89, \dots$$

is defined by setting the first two terms to 1 and making each subsequent term be the sum of the two preceding terms. Computations with this sequence turn

out to be a little simpler if we consider the sequence as starting with a “0th” term 0. Denoting the n^{th} term of the sequence by F_n , we have

n	F_n
0	0
1	1
2	1
3	2
4	3
5	5
6	8
7	13
8	21
9	34
10	55
11	89
⋮	

This is what is called a *recursive* definition—each term of the sequence is defined in terms of previous ones. But suppose we want to know what the 10574th term of the sequence is, without computing all the previous ones? It would be nice to have a formula that gives us F_n directly in terms of n .

There actually is such a formula; it is

$$F_n = \frac{1}{\sqrt{5}} \left(\frac{1 + \sqrt{5}}{2} \right)^n - \frac{1}{\sqrt{5}} \left(\frac{1 - \sqrt{5}}{2} \right)^n$$

The trouble with this formula is that it seems too cute; when you see it for the first time it is not at all clear where it could possibly come from, although you can verify for the first few values of n that it really does give the corresponding values of F_n correctly. So we will now explain where this formula comes from:

Consider the set V of *all* sequences $\{f_0, f_1, \dots\}$ of real numbers that satisfy

$$f_n = f_{n-1} + f_{n-2}$$

for all $n \geq 2$. For example, V contains the following sequences:

- 0, 1, 1, 2, 3, 5, 8, ...
- -1, 2, 1, 3, 4, 7, 11, ...
- 1, 0, 1, 1, 2, 3, 5, ...

Now V is actually a vector space. For,

- If we define the sum of two sequences $\{s_n\}$ and $\{t_n\}$ as the sequence $\{s_n + t_n\}$, then the sum of any two sequences in V is also a sequence in V .
- If c is a real number and we define $c\{s_n\}$ to be the sequence $\{cs_n\}$, then any real number times a sequence in V yields another sequence in V .

and this shows that V is a vector space over the field \mathbf{R} of real numbers. In fact, it is a 2-dimensional vector space. This is because each sequence in V is determined by its first two elements. If we make each sequence in V correspond to the ordered pair consisting of its first two elements, then we have an isomorphism between V and \mathbf{R}^2 . That is, the ordered pair $(0, 1)$ corresponds to the sequence

$$0, 1, 1, 2, 3, 5, \dots$$

and the ordered pair $(2, 1)$ corresponds to the sequence

$$2, 1, 3, 4, 7, 11, \dots$$

and so on. Addition of ordered pairs amounts to addition of the corresponding sequences, and multiplication by a scalar c acts similarly.

Now here's the idea: maybe we can find particular sequences in V that have a simple formula. In fact, we can. Suppose we look for a sequence that has the formula $f_n = t^n$ where t is some number to be determined. What must t satisfy? Well, since $f_{n+2} = f_{n+1} + f_n$ for all $n \geq 0$, we must have

$$t^{n+2} = t^{n+1} + t^n$$

for all $n \geq 0$. Factoring out t^n from each term in this equation, we see that we only need to find a t such that

$$t^2 = t + 1$$

There are two solutions to this quadratic equation:

$$t = \frac{1 \pm \sqrt{5}}{2}$$

and thus we have two specific sequences in V given by simple formulas:

$$1, \frac{1 + \sqrt{5}}{2}, \left(\frac{1 + \sqrt{5}}{2}\right)^2, \left(\frac{1 + \sqrt{5}}{2}\right)^3, \dots$$

and

$$1, \frac{1 - \sqrt{5}}{2}, \left(\frac{1 - \sqrt{5}}{2}\right)^2, \left(\frac{1 - \sqrt{5}}{2}\right)^3, \dots$$

Neither of the ordered pairs

$$\left(1, \frac{1 + \sqrt{5}}{2}\right) \quad \text{and} \quad \left(1, \frac{1 - \sqrt{5}}{2}\right)$$

is a multiple of the other. Therefore, these two ordered pairs constitute a basis of \mathbf{R}^2 , and any ordered pair can be constructed as a linear combination of these two ordered pairs. This amounts to saying that any sequence in V can be represented as a linear combination of these two sequences. To get the coefficients of the linear combination correct, all we have to do is check the first two terms of the sequence. For the Fibonacci sequence, we need to find a and b so that

$$(0, 1) = a \left(1, \frac{1 + \sqrt{5}}{2}\right) + b \left(1, \frac{1 - \sqrt{5}}{2}\right)$$

This amounts to picking a and b so that

$$\begin{aligned} a + b &= 0 & (n = 0) \\ a \left(\frac{1 + \sqrt{5}}{2}\right) + b \left(\frac{1 - \sqrt{5}}{2}\right) &= 1 & (n = 1) \end{aligned}$$

We can solve this to get

$$a = \frac{1}{\sqrt{5}} \quad b = -\frac{1}{\sqrt{5}}$$

Thus, the Fibonacci sequence has the representation

$$a \left(\frac{1 + \sqrt{5}}{2}\right)^n + b \left(\frac{1 - \sqrt{5}}{2}\right)^n$$

and substituting in the values of a and b we just found, this is just the formula we wrote previously.

Incidentally, this formula can be proved to be true by mathematical induction, but it is certainly clear that no one could possibly discover it that way.

1.7 Polynomials over a field

A polynomial (in one variable) over a field K is a finite sum of the form

$$a_0 + a_1x + a_2x^2 + a_3x^3 + \cdots + a_nx^n$$

where each coefficient a_i is an element of the field K . The *degree* of the polynomial is the highest power of the variable x that occurs in the sum; in this case, it is n .

Since we add polynomials by adding corresponding coefficients, and since we can multiply a polynomial by an element a of the field K by multiplying each coefficient by a , it is easy to see that the set of polynomials of degree $\leq n$ (for some fixed n) is a vector space over K . In fact, it is isomorphic to K^{n+1} (that is, it is just a renaming of K^{n+1}), since each polynomial can be equally well represented by the ordered $(n+1)$ -tuple of its coefficients

$$(a_0, a_1, a_2, \dots, a_n)$$

This is an $(n+1)$ -dimensional vector space over K .

For example, we see that the set of polynomials of degree 5 or less over \mathbf{Z}_2 is a vector space of degree 6 over \mathbf{Z}_2 , and therefore has 2^6 elements—there are 64 such polynomials. The first few are:

$$\begin{array}{l} 0 \\ 1 \\ \quad x \\ 1 + x \\ \quad \quad x^2 \\ 1 \quad + \quad x^2 \\ \quad \quad x + x^2 \\ 1 + x + x^2 \\ \quad \vdots \end{array}$$

The set of *all* polynomials over K (without any restriction on their degrees) is an infinite-dimensional vector space over K , since it is isomorphic to the set of tuples with infinitely many components

$$(a_0, a_1, a_2, \dots)$$

only finitely many of which are non-zero. This set of all polynomials in one variable x over the field K is denoted by $K[x]$.

Now $K[x]$ not only has the structure of a vector space—that is to say, we not only can add polynomials and multiply them by elements of the scalar field K —we can also multiply them by each other. That is, $K[x]$ is also a ring (in fact, an integral domain).

The fact that $K[x]$ is both a vector space and a ring is nothing to be confused about. We have already seen, for instance, that \mathbf{Z}_p is an additive group, but it is also a field. When we talk about $K[x]$ as a vector space, we are focussing our attention on its additive structure, on the operation of multiplying a scalar by a

polynomial, and on our ability to represent it as a set of tuples. When we talk about $K[x]$ as a ring, we are not concerned with multiplying by scalars per se, but we are very much concerned with multiplying two polynomials together.

Thinking of a set of polynomials as a vector space is particularly useful when we are concerned with dimensionality, or equivalently, with counting the number of polynomials there are. Thinking of a set of polynomials as a ring is particularly useful when we are concerned with questions of divisibility. Switching back and forth between these two points of view turns out to be quite productive.

1.8 Factoring polynomials

1.8.1 Division of polynomials

When we studied algebra in high school, we learned how to divide polynomials. Division of polynomials works just the same when the polynomials have coefficients in any field. For instance, Figure 1.3 shows a long division of two polynomials with coefficients in \mathbf{Z}_5 :

$$\begin{array}{r}
 4x^3 + 2x^2 + 2x + 1 \\
 3x^2 + 1 \overline{) 2x^5 + x^4 } \\
 \underline{2x^5 } \\
 x^4 + x^3 \\
 \underline{x^4 } \\
 x^3 + 3x^2 + 4x \\
 \underline{x^3 } \\
 3x^2 + 2x + 3 \\
 \underline{3x^2 } \\
 2x + 2
 \end{array}$$

Figure 1.3: Long division of two polynomials in \mathbf{Z}_5 .

The important fact is this: given two polynomials $P_1(x)$ and $P_2(x)$, P_1 can be divided by P_2 to get a quotient polynomial $Q(x)$ and a remainder polynomial $R(x)$ such that the degree of $R(x)$ is strictly less than that of $P_2(x)$, and we

have

$$P_1(x) = Q(x)P_2(x) + R(x)$$

For example, the long division computation above shows that in $\mathbf{Z}_5[x]$,

$$2x^5 + x^4 + 4x + 3 = (4x^3 + 2x^2 + 2x + 1)(3x^2 + 1) + (2x + 2)$$

and of course the degree of the remainder $2x + 2$ is strictly less than the degree of $3x^2 + 1$.

We use the notation $\deg(f)$ to denote the degree of the polynomial f . Thus, $\deg(3x^2 + 1) = 2$.

We can see now that division of polynomials acts very similarly to division of integers. Where in integers we would write

$$a = qb + r$$

with $0 \leq r < b$, here we have

$$\deg R < \deg P_2$$

This turns out to be enough to allow all the usual constructions to go through:

- The greatest common factor of two polynomials can be found by the Euclidean algorithm, just as it can for integers.
- Polynomials can be uniquely factored into polynomials that cannot be further factored. These polynomials are called *prime*, or *irreducible*, polynomials.

For example, in high school algebra, one might discover that

$$x^3 - x^2 + x - 1 = (x^2 + 1)(x - 1)$$

and neither polynomial on the right can be factored further. They are both prime polynomials.

Well, actually, that's a bit sloppy. The polynomial $x^2 + 1$ is a prime polynomial because in high school algebra we only deal with the real numbers. If we allow ourselves to use complex numbers, then $x^2 + 1$ is not prime. It factors like this:

$$x^2 + 1 = (x + i)(x - i)$$

In fact, that was the original reason complex numbers were introduced—to allow polynomials like $x^2 + 1$ to be factored.

Thus, *whether a polynomial is prime or not depends on what field we consider it as being over*. As an element of $\mathbf{R}[x]$, $x^2 + 1$ is prime. As an element of $\mathbf{C}[x]$, it is not.

1.8.2 Linear factors of polynomials

Thinking of a polynomial as a function of x (i.e., not just as a formal expression in the symbol x) can often help find linear factors of the polynomial. This is for the following simple reason:

The first-degree polynomial $p(x) = x - 7$ has the value 0 when $x = 7$; i.e., $p(7) = 7 - 7 = 0$. And more generally, if $p(x) = x - a$, then $p(a) = 0$.

We can push this a little farther: when a polynomial is factored into linear factors, we can tell at once what values of the variable x cause the polynomial to evaluate to 0. For instance, since (over the real numbers)

$$x^2 - x - 6 = (x + 2)(x - 3)$$

we can see at once that the function $p(x) = x^2 - x - 6$ satisfies

$$p(-2) = 0$$

$$p(3) = 0$$

We say the *zeros* of the polynomial p are -2 and 3 .

Well, that's not too impressive, really. But suppose we have figured out that the polynomial

$$p(x) = x^4 - 2x^3 - 41x^2 + 42x + 360$$

can be factored as

$$p(x) = (x + 3)(x - 4)(x + 5)(x - 6)$$

Then we know at once that the zeros of p are -3 , 4 , -5 , and 6 :

$$p(-3) = 0$$

$$p(4) = 0$$

$$p(-5) = 0$$

$$p(6) = 0$$

This works in reverse, too: if we have been given the polynomial

$$p(x) = x^4 - 2x^3 - 41x^2 + 42x + 360$$

and we know it vanishes when x is -3 , 4 , -5 , and 6 , then we know at once that p factors as

$$p(x) = (x + 3)(x - 4)(x + 5)(x - 6)$$

This is very useful as a way of finding the factors of a polynomial, because in many cases it is not hard to find the zeros of the polynomial.

We can state this result as a couple of simple but important theorems:

1.15 Theorem (Remainder theorem) *If $p(x)$ is a polynomial, then the remainder when $p(x)$ is divided by the polynomial $x - a$ is just $p(a)$.*

PROOF. The remainder will be a constant, because it is a polynomial whose degree is less than that of $x - a$. Thus, we have

$$p(x) = q(x)(x - a) + r$$

Substituting a for x , we get

$$p(a) = q(a)(a - a) + r = r$$

□

1.16 Theorem (Factor theorem) *If $p(x)$ is a polynomial, then $x - a$ is a factor of $p(x)$ if and only if $p(a) = 0$.*

PROOF. Using the notation of the previous proof, $x - a$ is a factor of $p(x)$ if and only if the remainder r is 0; i.e., if and only if $p(a) = 0$. □

Nothing in these proofs uses any special properties of the field over which the polynomial is defined. Thus, these results are true over any field. For example, over the field \mathbf{Z}_2 , the polynomial $x^2 + 1$ vanishes for $x = 1$. This is because all computations now are (mod 2); we are just saying that

$$1^2 + 1 \equiv 0 \pmod{2}$$

Therefore $x - 1$ (which is the same as $x + 1$ over \mathbf{Z}_2 , because -1 is the same as 1 in this field) is a factor of $x^2 + 1$; in fact, we have

$$x^2 + 1 = (x + 1)(x + 1)$$

in $\mathbf{Z}_2[x]$, since

$$(x + 1)(x + 1) = x^2 + 2x + 1$$

and the middle term on the right vanishes since it is a multiple of 2.

Thus, $x^2 + 1$ factors into linear factors in $\mathbf{Z}_2[x]$. Note, however, that in $\mathbf{R}[x]$, the polynomial $x^2 + 1$ is prime, as we have already mentioned.

The fact that $x^2 + 1 = (x + 1)(x + 1)$ in $\mathbf{Z}_2[x]$ is just a rewording of the result already noted in Theorem 1.5 (page 11). In fact, for any prime p and any $n > 0$ we have

$$(x + a)^{p^n} = x^{p^n} + a^{p^n}$$

in $\mathbf{Z}_p[x]$.

The factor theorem has the following consequence:

1.17 Theorem *A polynomial of degree n has at most n roots.*

PROOF. This is because each root a corresponds to a factor of the form $x - a$. If there were more than n such factors, then multiplying them together would yield a polynomial of degree greater than n . This polynomial would either be the original polynomial, or it would be a factor of it. Either way, the original polynomial would have degree greater than n , a contradiction. \square

As a remarkable consequence of this simple result, we have

1.18 Theorem *The multiplicative group \mathbf{Z}_p^* is cyclic.*

PROOF. The reason this is true is that \mathbf{Z}_p^* is the multiplicative group of the field \mathbf{Z}_p . Let us look at the orders of the elements of \mathbf{Z}_p^* . Say n is the largest order of any element in \mathbf{Z}_p^* , and say a is such an element (i.e., the order of a is n). Corollary 1.12 (page 22) shows that every element of \mathbf{Z}_p^* has order dividing n . Therefore, every element b of \mathbf{Z}_p^* satisfies $b^n = 1$. To put it another way, the polynomial $x^n - 1$ has $p - 1$ roots in \mathbf{Z}_p^* . This can only be true if $n \geq p - 1$. On the other hand, the order of any element divides the order of \mathbf{Z}_p^* , so actually $n = p - 1$. Therefore, the order of a is exactly $p - 1$, and this element therefore generates all of \mathbf{Z}_p^* , so \mathbf{Z}_p^* is cyclic. \square

An element a of \mathbf{Z}_p^* whose order is $p - 1$ (and which is therefore a generator of \mathbf{Z}_p^*) is called a *primitive root* (or a *primitive element*) of the field \mathbf{Z}_p , or alternatively, a primitive root (mod p). (The term “root” comes from the fact that it is a root of the polynomial $x^{p-1} - 1$.)

We have just shown that each field \mathbf{Z}_p has at least one primitive root. Actually, there are a number of them, and it is known how many. But there is no simple way to find them. Of course, trial and error always works. As an example of the kind of behavior we can expect, let us consider the field \mathbf{Z}_7 . This field has 6 non-zero elements. Since its multiplicative group has 6 elements and is cyclic, it is isomorphic to the additive group \mathbf{Z}_6 . (Actually, every abelian group with 6 elements is isomorphic to \mathbf{Z}_6 , but that kind of thing is not true in general—remember that $\mathbf{Z}_2 \times \mathbf{Z}_2$ is a group with 4 elements that is not isomorphic to \mathbf{Z}_4 .) We have already seen that \mathbf{Z}_6 has two elements of order 6. Therefore, \mathbf{Z}_7^* similarly has two elements of (multiplicative) order 6. They are 3 and 5: we have (remember that everything here is (mod 7))

n	0	1	2	3	4	5
3^n	1	3	2	6	4	5
5^n	1	5	4	6	2	3

We can use exactly the same reasoning as in the proof of the previous theorem to show the following generalization, which is the key result for our purposes here:

1.19 Theorem *The multiplicative group of any finite field is cyclic.*

The proof is almost word-for-word the same. It just depends on the facts that

- A polynomial of degree n over a field cannot have more than n roots.
- In a finite abelian group, if a is an element whose order is maximal, the order of any other element divides the order of a .

In the next chapter we will show how to construct many new finite fields. But this result does not depend on any construction, so we have put it here.

1.8.3 Formal derivatives

If $p(x)$ is a polynomial, we say that a is a *simple root* or *simple zero* of p if $(x-a)|p(x)$ (so certainly $p(a) = 0$), but $(x-a)^2$ does not divide $p(x)$. Otherwise, a is a *multiple root* of $p(x)$. It turns out that there is an easy way to tell if a root is simple or not. To understand this, let us think of polynomials over \mathbf{R} as functions. If $(x-a)^2$ divides $p(x)$, then near a , $p(x)$ is tangent to the x -axis; that is, its derivative is 0 at a . On the other hand, if a is a simple root, then the derivative will not be 0; $p(x)$ will pass through the x -axis at a with a positive or negative slope.

It turns out that, even when we consider polynomials (mod p), we can use this same technique. First we have to say what we mean by the derivative of a polynomial: If

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

is a polynomial over a field K , then the derivative of p is defined to be the polynomial

$$p'(x) = n a_n x^{n-1} + (n-1) a_{n-1} x^{n-2} + \cdots + a_1$$

It's important to realize that this is a purely algebraic notion. There is no question of the derivative being the slope of anything, or being a limit of a difference quotient, as it is in calculus. (In fact, we really have no notion of limit in an abstract field.) Nevertheless, the derivative is quite useful.

The first thing to know about this definition of the derivative is that it satisfies all the usual formulas we are familiar with. For instance,

$$\begin{aligned} ((p(x) + q(x))' &= p'(x) + q'(x) \\ (p(x)q(x))' &= p(x)q'(x) + p'(x)q(x) \end{aligned}$$

We won't bother to prove these results, but they are easy to prove using induction. (This is much more elementary than in calculus, where these formulas are proved for arbitrary differentiable functions—here we only need to differentiate polynomials, which are much simpler.)

So the derivative of a polynomial is easy to compute. Here is how it is used:

1.20 Theorem *A root a of the polynomial $p(x)$ is a multiple root if and only if $p'(a) = 0$.*

PROOF. Since a is a root of $p(x)$, we have $p(x) = (x - a)g(x)$ where $g(x)$ is a polynomial. a is a multiple root of p if and only if it is a root of g .

Now taking the derivative, we have

$$p'(x) = (x - a)g'(x) + g(x)$$

When x is a , this becomes $p'(a) = g(a)$. Thus, a is a root of g if and only if $p'(a) = 0$. \square

Chapter 2

Finite Fields

2.1 Irreducible polynomials and minimal polynomials

When we are talking about polynomials over a field, most of what we care about does not change if the polynomial is multiplied by a constant. For instance, if $f(x)$ is a polynomial over K and α is a root of $f(x)$, then α is also a root of $cf(x)$ for any element c of K . Similarly, if $f(x)$ is irreducible, then so is $cf(x)$. Therefore, we often normalize polynomials by dividing by the coefficient of their highest-degree term. After this normalization, the coefficient of the highest-degree term is 1. Such polynomials are called *monic* polynomials. For instance,

$$3x^2 - 4$$

is not a monic polynomial over \mathbf{Z}_5 . But after dividing (in \mathbf{Z}_5) by the coefficient 3 of x^2 , we get the monic polynomial

$$x^2 - 3$$

Suppose that K and L are two fields, and that K is a subfield of L (so they have the same elements 0 and 1, and the same addition and multiplication). Suppose that p is a monic polynomial over K (equivalently, p is a monic polynomial in $K[x]$). Suppose further that p is irreducible over K . We can show that if p has a root in L —say $p(\alpha) = 0$ for some α in L —then p is the unique monic polynomial of smallest degree having this property.

For let f be a monic polynomial over K of minimal degree having α as a root.

Dividing p by f , we have

$$p(x) = q(x)f(x) + r(x)$$

where $\deg(r) < \deg(f)$. Substituting α for x , we see that $r(\alpha) = 0$. But by the assumption that f was a monic polynomial of smallest degree having α as a root, we must have $r = 0$. This means that f divides p . But p is irreducible. Therefore, f must be p . (Up to a constant multiple; but since both f and p are both monic, they are in fact identical.)

Thus, a monic irreducible polynomial is the monic polynomial of minimal degree (or the *minimal polynomial*, as we say) of each of its roots.

Conversely, a minimal polynomial must be irreducible. That is, suppose that f is the monic polynomial of smallest degree having α as a root. If f is not irreducible, then it factors into two other monic polynomials: $f(x) = g(x)h(x)$. Substituting α for x , we have $0 = f(\alpha) = g(\alpha)h(\alpha)$, and so at least one of $g(\alpha)$ and $h(\alpha)$ must be 0; but this contradicts the minimality of the degree of f .

Thus, we have shown:

2.1 Theorem *A monic polynomial over K is irreducible if and only if it is the minimal polynomial of a root in an extension field L of K .*

Note, by the way, that in the case of polynomials of degree 1, the extension field L could just be K itself. For instance, the polynomial $x - 1$ is irreducible over any field, and it is clearly the minimal polynomial of its root 1 (which is in every field). The interesting cases of this, however, are those in which the polynomial has degree greater than 1 and in which L properly contains K .

2.2 Corollary *If $p(x)$ is an irreducible polynomial over a field K , and if α is a root of p in an extension field of K , and if $f(x)$ is another polynomial over K such that $f(\alpha) = 0$, then $p(x) \mid f(x)$.*

PROOF. Divide f by p ; we get

$$f(x) = q(x)p(x) + r(x)$$

where $\deg(r) < \deg(p)$. Substituting α for x , we find that $r(\alpha) = 0$. Since p is the minimal polynomial for α , we must have $r = 0$, so $p(x) \mid f(x)$. \square

2.2 Algebraic field extensions

Now we are going to show how to make new fields from old ones—to be precise, how to extend fields to make larger ones. The way to think about this is to remember how the complex number field \mathbf{C} is constructed from the real field \mathbf{R} .

We start with the polynomial $p(x) = x^2 + 1$ which is irreducible over \mathbf{R} . (An equivalent way of expressing this is to say that it is irreducible in $\mathbf{R}[x]$.) We then introduce a new number, which we denote by the symbol i , and which is such that $i^2 = -1$. That is, $i^2 + 1 = 0$, or to put it another way, i is a root of $p(x)$. We then consider the set of all polynomials in i . Actually, any polynomial in i has degree 1: it looks like $a + bi$, because any higher power of i can be reduced by successive application of the identity $i^2 = -1$. For instance,

$$\begin{aligned} 2 - 3i + 4i^2 - 5i^3 + 6i^4 &= 2 - 3i + 4(-1) - 5i(-1) + 6(-1)^2 \\ &= 4 + 2i \end{aligned}$$

Another way to look at this is to reduce polynomials in i by adding or subtracting multiples of $i^2 + 1$ (which of course has been defined to be 0). So for instance, using the same example:

$$\begin{aligned} 2 - 3i + 4i^2 - 5i^3 + 6i^4 &= 2 - 3i \\ &\quad + 4i^2 - 4(i^2 + 1) \\ &\quad - 5i^3 + 5i(i^2 + 1) \\ &\quad + 6i^4 - 6i^2(i^2 + 1) + 6(i^2 + 1) \\ &= 4 + 2i \end{aligned}$$

Now we are so used to thinking of the complex number i as being an honest number—an actual point in the complex plane—that we allow ourselves to forget that originally, it was just a symbol. (In fact, historically, the complex plane came later.) We could just as well have used any symbol. We could just as well have used the symbol x . This would make it more obvious that what we are doing is just reducing polynomials in x by adding and subtracting multiples of the polynomial $x^2 + 1$:

$$\begin{aligned} 2 - 3x + 4x^2 - 5x^3 + 6x^4 &= 2 - 3x \\ &\quad + 4x^2 - 4(x^2 + 1) \\ &\quad - 5x^3 + 5x(x^2 + 1) \\ &\quad + 6x^4 - 6x^2(x^2 + 1) + 6(x^2 + 1) \\ &= 4 + 2x \end{aligned}$$

Once we have gotten this far, it's only a small step to see that what we are “really” doing is computing the remainder when dividing by the polynomial $x^2 + 1$:

$$2 - 3x + 4x^2 - 5x^3 + 6x^4 = (6x^2 - 5x - 2)(x^2 + 1) + (4 + 2x)$$

That is, we are throwing away the multiple $(6x^2 - 5x - 2)(x^2 + 1)$ of $x^2 + 1$, leaving the remainder $4 + 2x$.

This is entirely analogous to forming \mathbf{Z}_p by taking congruence classes in the ring \mathbf{Z} with respect to the modulus p (in that case, we throw away multiples of the prime number p); here we are taking congruence classes in the ring $\mathbf{R}[x]$ with respect to the modulus $x^2 + 1$.

We can perform this construction in any polynomial ring: if K is a field, and if $p(x)$ is an irreducible polynomial in $K[x]$, we can form the congruence classes with respect to the modulus $p(x)$. The way algebraists usually write this set of congruence classes is like this:

$$K[x]/(p(x)) \quad \text{or} \quad K[x]/p(x)K[x]$$

the idea being that $p(x)K[x]$ just means the set of all multiples of $p(x)$ in $K[x]$, that is,

$$p(x)K[x] = \{p(x)q(x) : q(x) \in K[x]\}$$

and this set of polynomials is the 0 congruence class in the family of congruence classes we are constructing. As far as algebraists are concerned, the set \mathbf{C} of complex numbers is just the set of congruence classes $\mathbf{R}[x]/(x^2 + 1)\mathbf{R}[x]$.

Incidentally, and for the same reason, another notation for what we have written as \mathbf{Z}_p is $\mathbf{Z}/p\mathbf{Z}$.

Now as we know, the set of congruence classes (mod p) is conventionally denoted $\{0, 1, \dots, p - 1\}$. There is really an abuse of notation here: the symbol “0”, which ordinarily represents the integer 0, is here being used to represent the congruence class $p\mathbf{Z}$. We can do this because the number 0 is the unique number n in that congruence class that satisfies $0 \leq n \leq p - 1$. And in general, the congruence class that an arbitrary number n belongs to (mod p) is just the congruence class that we represent by the number which is the remainder when n is divided by p (the remainder r being in the interval $0 \leq r < p$).

In the same way, the set of congruence classes $K[x]/p(x)K[x]$ can be represented by the set of polynomials of degree less than n , where n is the degree of the polynomial $p(x)$. Any polynomial f belongs to the congruence class represented by its remainder when divided by $p(x)$.

In this case, however, this abuse of language can be a bit confusing—when we see the symbol “ x ”, we normally do not think of a congruence class, let alone a root of a polynomial. So what we conventionally do is introduce a new symbol. It could be any symbol; let us use κ . Then instead of writing a polynomial $f(x)$ with the understanding that this $f(x)$ is congruent to a polynomial of degree less than n , we write $f(\kappa)$, with the understanding that $f(\kappa)$ equals a polynomial in κ of degree less than n . This is entirely analogous to what goes on when we

pass from saying

$$12 \equiv 5 \pmod{7}$$

to

$$12 = 5 \quad \text{in } \mathbf{Z}_7$$

In this formulation, $p(\kappa)$ actually equals 0, and so κ is a root of the polynomial $p(x)$. This process is called *adjoining* a root κ of $p(x)$ to the field K (just as \mathbf{C} is produced by adjoining the root i of $x^2 + 1$ to \mathbf{R}).

Now notice:

- $K[\kappa]$ consists of all the polynomials of degree less than n over K .
- Adding elements of $K[\kappa]$ just consists of adding corresponding coefficients of the polynomials
- Multiplying an element of $K[\kappa]$ by an element λ of K can be performed by multiplying each coefficient of the polynomial in $K[\kappa]$ by λ .

Therefore, $K[\kappa]$ is a vector space over K , of dimension n . In particular, if K is a finite field, then $K[\kappa]$ is also finite.

$K[\kappa]$ is also a ring—it is closed under multiplication. To multiply two elements of $K[\kappa]$, we just multiply the polynomials as usual and then reduce by dividing by the polynomial p . That is, we replace the product by the remainder we get by dividing it by the polynomial $p(x)$.

To be honest about this, we should note that addition and multiplication of congruence classes of polynomials mod $p(x)$ are both well-defined. This just amounts to proving analogues of Theorems 1.1 and 1.2. The proofs are exactly the same.

Now *provided* p is a prime, $\mathbf{Z}_p = \mathbf{Z}/p\mathbf{Z}$ is not only a ring—it is actually a field. We saw this back in Chapter 1. The same holds true for $K[x]/p(x)K[x]$:

2.3 Theorem *If K is a finite field, and if $p(x)$ is an irreducible polynomial over K , then $K[x]/p(x)K[x]$ (which we denote by $K[\kappa]$) is a field.*

Remark This theorem is really true for *any* field K . But the proof when K is a finite field is much more elementary than the proof in general, and it is all we need here.

PROOF. $K[\kappa]$ is obviously a ring—it is closed under multiplication, and so on. To show it is a field, we only need to prove that every non-zero element in $K[\kappa]$ has a multiplicative inverse.

The proof is essentially the same as that which shows that \mathbf{Z}_p is a field:

$K[\kappa]$ is a finite set, as we have just seen. Say its non-zero elements are enumerated as $\{f_1(\kappa), f_2(\kappa), \dots, f_m(\kappa)\}$ (for some finite number m). That is, this set is just the set of non-zero polynomials over K of degree less than n . Now if g is any element of this set, consider the set of polynomials that we get by multiplying each of these polynomials by g (and then reducing by $p(x)$ as usual). This set is

$$\{g(\kappa)f_1(\kappa), g(\kappa)f_2(\kappa), \dots, g(\kappa)f_m(\kappa)\}$$

No two elements of this set are the same. For if $g(\kappa)f_i(\kappa) = g(\kappa)f_j(\kappa)$, for instance, then we would have

$$g(\kappa)(f_i(\kappa) - f_j(\kappa)) = 0$$

or equivalently,

$$p(x) | g(x)(f_i(x) - f_j(x))$$

But since $p(x)$ is irreducible and g is a non-zero polynomial of degree less than the degree of p , this means that $p(x)$ divides $f_i(x) - f_j(x)$, and since $f_i(x) - f_j(x)$ also has degree less than that of $p(x)$, this means that i and j must be equal.

Thus, multiplication by $g(\kappa)$ just permutes the finite set of non-zero elements of $K[\kappa]$, and hence there is some i for which $g(\kappa)f_i(\kappa) = 1$. This shows that $g(\kappa)$ has a multiplicative inverse in $K[\kappa]$, and hence $K[\kappa]$ is a field. \square

Thus if $p(x)$ is irreducible over K , $K[x]/p(x)K[x] = K[\kappa]$ is a field that contains K as a subfield. It is called an *algebraic field extension* of K .

As we mentioned, this proof does not work for infinite fields K —for instance, this proof cannot be used to show that \mathbf{C} , constructed from \mathbf{R} , is a field. But it works for all the fields we care about here.

κ was constructed so that it is a root of the irreducible polynomial $p(x)$ of degree n . The only polynomials that have κ as a root are the polynomials that are divisible by $p(x)$ —any other polynomial $f(x)$ has a non-zero remainder when divided by $p(x)$, and by definition $f(\kappa)$ is represented by that remainder and is therefore non-zero.

Here is an example of this construction: Let us start with \mathbf{Z}_2 , the simplest finite field. Over this field, the polynomial $p(x) = x^2 + x + 1$ is irreducible. (In this case, one can verify this fact by simply trying to divide it by polynomials of smaller degree.) Adjoining a root κ of this polynomial to \mathbf{Z}_2 yields a field containing four elements: 0, 1, κ , and $\kappa + 1$. (These are the four polynomials in κ of degree 0 or 1 over \mathbf{Z}_2 .) These elements add and multiply as follows:

+	0	1	κ	$\kappa + 1$
0	0	1	κ	$\kappa + 1$
1	1	0	$\kappa + 1$	κ
κ	κ	$\kappa + 1$	0	1
$\kappa + 1$	$\kappa + 1$	κ	1	0

\cdot	0	1	κ	$\kappa + 1$
0	0	0	0	0
1	0	1	κ	$\kappa + 1$
κ	0	κ	$\kappa + 1$	1
$\kappa + 1$	0	$\kappa + 1$	1	κ

You can see just by looking at these tables that $\mathbf{Z}_2[\kappa]$ is really a field. Note that its additive group (i.e., just considering it as a group under addition) is isomorphic to $\mathbf{Z}_2 \times \mathbf{Z}_2$. This is because it is the same as the polynomials of degree 1 over \mathbf{Z}_2 , and these polynomials add by adding the corresponding coefficients. Equivalently, we just have a vector space of dimension 2 over \mathbf{Z}_2 . The multiplicative group of this field has three elements (1, κ , and $\kappa + 1$), and is cyclic—it is isomorphic to \mathbf{Z}_3 .

Also note that since addition and multiplication are both commutative, these tables are both symmetric about the main diagonal. So we really only have to fill in the elements on and above that diagonal to specify these operations.

We want to emphasize that it is crucial in this construction that the polynomial $p(x)$ be irreducible over K . For if it is not, say $p(x) = f(x)g(x)$ where f and g are each non-zero polynomials whose degrees are less than the degree of p . Then we can still adjoin a root κ of $p(x)$. But we will have $f(\kappa)g(\kappa) = p(\kappa) = 0$, so $f(\kappa)$ and $g(\kappa)$ are non-zero elements in $K[\kappa]$ whose product is 0, which shows that $K[\kappa]$ cannot possibly be a field.

For example, consider the polynomial $f(x) = x^2 + 2$ over \mathbf{Z}_3 . We can adjoin a root κ to form $\mathbf{Z}_3[\kappa]$. The addition and multiplication tables in $\mathbf{Z}_3[\kappa]$ look like this:

+	0	1	2	κ	$\kappa + 1$	$\kappa + 2$	2κ	$2\kappa + 1$	$2\kappa + 2$
0	0	1	2	κ	$\kappa + 1$	$\kappa + 2$	2κ	$2\kappa + 1$	$2\kappa + 2$
1		2	0	$\kappa + 1$	$\kappa + 2$	κ	$2\kappa + 1$	$2\kappa + 2$	2κ
2			1	$\kappa + 2$	κ	$\kappa + 1$	$2\kappa + 2$	2κ	$2\kappa + 1$
κ				2κ	$2\kappa + 1$	$2\kappa + 2$	0	1	2
$\kappa + 1$					$2\kappa + 2$	2κ	1	2	0
$\kappa + 2$						$2\kappa + 1$	2	0	1
2κ							κ	$\kappa + 1$	$\kappa + 2$
$2\kappa + 1$								$\kappa + 2$	κ
$2\kappa + 2$									$\kappa + 1$

·	0	1	2	κ	$\kappa + 1$	$\kappa + 2$	2κ	$2\kappa + 1$	$2\kappa + 2$
0	0	0	0	0	0	0	0	0	0
1		1	2	κ	$\kappa + 1$	$\kappa + 2$	2κ	$2\kappa + 1$	$2\kappa + 2$
2			1	2κ	$2\kappa + 2$	$2\kappa + 1$	κ	$\kappa + 2$	$\kappa + 1$
κ				1	$\kappa + 1$	$2\kappa + 1$	2	$\kappa + 2$	$2\kappa + 2$
$\kappa + 1$					$2\kappa + 2$	0	$2\kappa + 2$	0	$\kappa + 1$
$\kappa + 2$						$\kappa + 2$	$\kappa + 2$	$2\kappa + 1$	0
2κ							1	$2\kappa + 1$	$\kappa + 1$
$2\kappa + 1$								$\kappa + 2$	0
$2\kappa + 2$									$2\kappa + 2$

We see that $\mathbf{Z}_3[\kappa]$ has divisors of zero, so it is certainly not a field. This is because the polynomial $x^2 + 2$ factors over \mathbf{Z}_3 :

$$x^2 + 2 = (x + 1)(x + 2)$$

Now given an irreducible polynomial $p(x)$, this construction which adjoins a root κ of $p(x)$ to form $K[\kappa]$ is quite general, in the following respect:

Suppose that

- K is a finite field containing d elements.
- $p(x)$ is an irreducible polynomial of degree n over K .
- L is a field containing K as a subfield (L is an *extension* of the field K).
- There is an element α of L that is a root of $p(x)$.

Let $K[\alpha]$ denote the set of all polynomials in α with coefficients in K . $K[\alpha]$ is a subset of the field L . It is clear that $K[\alpha]$ is closed under addition, subtraction, and multiplication, but it is not clear that non-zero elements of $K[\alpha]$ have

multiplicative inverses in $K[\alpha]$. (They do of course have multiplicative inverses in the field L , by assumption.) We will show this is true, but for now we will not assume it. In any case, $K[\alpha]$ has *no more than* d^n elements, because each element of $K[\alpha]$ has a representation as a polynomial over K of degree less than n , and there are d^n such polynomials. (We will see in a moment that $K[\alpha]$ actually has exactly d^n elements.)

Now let $K[\kappa]$ be the field we constructed above. We know that $K[\kappa]$ has exactly d^n elements. We will show that $K[\kappa]$ is isomorphic to $K[\alpha]$, the isomorphism being such that κ is identified with α .

The proof of this goes as follows:

1. Each element of $K[\kappa]$ is a polynomial $f(\kappa)$. To each such element, we associate the element $f(\alpha)$ of $K[\alpha]$. Then addition and multiplication in $K[\kappa]$ correspond to addition and multiplication in $K[\alpha]$ of these elements (since they just correspond to addition and multiplication of corresponding polynomials, with reduction by p after each multiplication). Thus $K[\kappa]$ is mapped into $K[\alpha]$ by this correspondence in a manner that preserves the operations of addition and multiplication.
2. The only element of $K[\kappa]$ which is identified with 0 in $K[\alpha]$ is 0 itself. For an element of $K[\kappa]$ is just a polynomial $f(\kappa)$ in κ . If there is any such non-zero polynomial f in $K[x]/p(x)K[x]$ such that $f(\alpha) = 0$ in $K[\alpha]$, let f be such a polynomial of smallest degree. Dividing $p(x)$ by $f(x)$, we get

$$p(x) = q(x)f(x) + r(x)$$

and setting $x = \alpha$, we see that $r(\alpha) = 0$. But r has degree less than the degree of f , and so since the degree of f was minimal, r must be the zero polynomial. But this in turn means that f divides p , which is a contradiction, because p is irreducible.

Actually, this was just a reworking of the proof of Section 2.1. We could just as well have argued as follows: If there is any non-zero polynomial f in $K[x]/p(x)K[x]$ such that $f(\alpha) = 0$ in $K[\alpha]$ (and hence in the field L), then by Corollary 2.2, p divides f ; that is $f(x) = p(x)q(x)$. Substituting κ for x , we see that $f(\kappa) = p(\kappa)q(\kappa) = 0$, so $f(\kappa)$ really was 0; or equivalently, f really was 0 in $K[x]/p(x)K[x]$.

3. Distinct elements of $K[\kappa]$ correspond under this correspondence with distinct elements of $K[\alpha]$. For if there are two elements of $K[\kappa]$ —call them $f(\kappa)$ and $g(\kappa)$ —such that $f(\alpha) = g(\alpha)$, then the polynomial $h = f - g$ satisfies $h(\alpha) = 0$. By what we just showed, h must be 0; i.e., $g = h$.
4. Therefore $K[\alpha]$ must have *no less than* d^n elements, and so it has exactly d^n elements and the correspondence we have just set up identifies all of

$K[\kappa]$ with all of $K[\alpha]$. So $K[\alpha]$ is isomorphic to $K[\kappa]$. In particular, it is a field—each non-zero element has a multiplicative inverse in $K[\alpha]$.

Thus, $K[\alpha]$ is actually a field extension of K and a subfield of L ; we have

$$K \subset K[\alpha] \subset L$$

Now this result is true for *every* root of the irreducible polynomial $p(x)$ in *every* extension field L of K —if α is any such root, $K[\alpha]$ is isomorphic to $K[\kappa]$, with α corresponding to κ . This result is really quite significant. For instance, it enables us to prove the following result immediately:

2.4 Theorem *If $p(x)$ is an irreducible polynomial over K , and if α is a root of $p(x)$ in some extension field, and if the multiplicative order of α is r (that is, if $\alpha^i = 1$ if and only if i is a multiple of r), then the multiplicative order of any root of $p(x)$ in any extension field of K is also r .*

PROOF. The multiplicative order of α has to be the same as the multiplicative order of κ in $K[\kappa]$. But this is true for any other root of $p(x)$ also. \square

2.3 Finite fields

As a result of all we have proved in the last section, we see that we can construct a finite field of characteristic p by finding an irreducible polynomial $p(x)$ over \mathbf{Z}_p and adjoining a root of this polynomial to \mathbf{Z}_p . Equivalently, given such a polynomial $p(x)$, of degree n , just consider the set of polynomials of degree less than n with coefficients in \mathbf{Z}_p , and perform multiplication (mod $p(x)$).

This then raises the question of how to find irreducible polynomials over \mathbf{Z}_p . In particular, is there guaranteed to be an irreducible polynomial of degree n for every positive integer n ?

It turns out that there is, but this fact is not obvious. Nevertheless, we can show that plenty of finite fields exist by an indirect method, which then will show that such irreducible polynomials really exist.

Let K be a finite field. We have already seen that K has p^n elements, where n is its dimension over its prime field and where p is its characteristic.

Since the multiplicative group K^* has order $p^n - 1$, every element x of K^* must satisfy $x^{p^n - 1} = 1$. (This is simply because the order of any element divides the order of the group K^* .) This means that each element a of K^* is a root of the polynomial $x^{p^n - 1} - 1$, and so $x - a \mid x^{p^n - 1} - 1$ in $\mathbf{Z}_p[x]$. There are exactly $p^n - 1$ elements a of K^* , and $x^{p^n - 1} - 1$ has each $x - a$ as a factor, so this forms the

complete factorization of this polynomial:

$$x^{p^n-1} - 1 = \prod_{a \in K^*} (x - a)$$

This already, by the way, has a famous consequence in elementary number theory. Take K to be simply the field \mathbf{Z}_p . In this case, $n = 1$, and the elements of K^* are just the numbers $\{1, 2, \dots, p-1\}$. We have

$$x^{p-1} - 1 = \prod_{i=1}^{p-1} (x - i)$$

Equating the constant terms on each side, this yields

2.5 Theorem (Wilson's Theorem) *If p is a prime,*

$$(p-1)! \equiv -1 \pmod{p}$$

(For instance, $4! = 24 \equiv -1 \pmod{5}$.)

Now let us reverse this process. To construct a finite field of dimension n over \mathbf{Z}_p , we need to adjoin roots of $x^{p^n-1} - 1$ to \mathbf{Z}_p . It is a little easier if we consider instead the polynomial $x^{p^n} - x$ —this just introduces a factor of x , and so now *every* element of K (including 0) will be a root of this polynomial. We will construct a field in which this polynomial factors into linear factors as follows:

First factor $x^{p^n} - x$ into irreducible factors over \mathbf{Z}_p . Pick one of those factors of degree greater than 1, and adjoin a root κ_1 to form a field $\mathbf{Z}_p[\kappa_1]$. Now factor $x^{p^n} - x$ into irreducible factors over this new field. It will have at least one more linear factor (namely, $x - \kappa_1$) than before. Pick an irreducible factor of degree greater than 1 and repeat the process. We build up successively larger and larger fields until we get—in a finite number of steps—a field L in which $x^{p^n} - x$ factors into linear factors. Say these factors are

$$x^{p^n} - x = \prod_{i=1}^{p^n} (x - a_i)$$

Now none of these factors are repeated; that is, no a_i occurs twice. This is because since the derivative of the polynomial $f(x) = x^{p^n} - x$ is $f'(x) = p^n x^{p^n-1} - 1 = -1$, $f(x)$ and $f'(x)$ have no common roots—in fact, $f'(x)$ is never 0 at all. Therefore there are p^n distinct elements a_i .

Let F be the set consisting all the a_i . F is a subset of the field L containing p^n elements. F certainly contains 0, since $0^{p^n} - 0 = 0$, and for the same reason, F contains 1. We show that F is a field:

- F is closed under multiplication, because if a and b are in F , we have $a^{p^n} = a$ and $b^{p^n} = b$, and so

$$(ab)^{p^n} = a^{p^n} b^{p^n} = ab$$

which shows that ab is also in F .

- F is closed under addition, because if a and b are in F , we have

$$(a + b)^{p^n} = a^{p^n} + b^{p^n}$$

(Expanding the left-hand side by the binomial theorem, all the coefficients of the “middle” terms are divisible by p , so those terms drop out in any field of characteristic p .)

- -1 is in F , because

$$(-1)^{p^n} = \begin{cases} -1 & \text{if } p \text{ is an odd prime} \\ 1 \equiv -1 \pmod{2} & \text{if } p = 2 \end{cases}$$

Therefore, if a is in F , also $-a = (-1)a$ is in F .

- If a is in F and $a \neq 0$, then a has a multiplicative inverse. This follows by the usual argument: the set of non-zero elements is a finite set closed under multiplication, and multiplication by any one of them permutes the set.

Thus, F is a field containing exactly p^n elements.

Now that we have proved that a finite field with p^n elements exists for each prime p and each positive integer n , we can show that for each such p and n , there is an irreducible polynomial $p(x)$ over \mathbf{Z}_p of order n , which can be used to construct such a field directly:

2.6 Theorem *For each prime p and positive integer n , there is an irreducible polynomial $p(x)$ in $\mathbf{Z}_p[x]$ of degree n .*

PROOF. Let K be a field over \mathbf{Z}_p with p^n elements, as constructed above. We know that this field has a primitive element. Call it a . We know that there is at least one non-zero polynomial over \mathbf{Z}_p having a as a root, namely, $x^{p^n} - x$. Therefore there is one of smallest degree. Call it $p(x)$, and say $p(x)$ has degree m . $p(x)$ must be irreducible, as we have already seen in Theorem 2.1.

Now the dimension of $\mathbf{Z}_p[a]$ over \mathbf{Z}_p is just m , since $\mathbf{Z}_p[a]$ just consists of all the polynomials in a of degree less than m .

On the other hand, $\mathbf{Z}_p[a]$ certainly includes all the powers of a . Since a is a primitive element, these are all the elements of K . So $\mathbf{Z}_p[a] = K$. (If a were

not a primitive element, $\mathbf{Z}_p[a]$ would still be a finite field, but it might not be all of K —it might just be a subfield of K .) Therefore, the dimension of $\mathbf{Z}_p[a]$ over \mathbf{Z}_p must be n .

This shows that $m = n$, i.e., that $p(x)$ is an irreducible polynomial over \mathbf{Z}_p of degree n . \square

Finally, we can show that any two finite fields of order p^n are isomorphic, and that the ones we have constructed are all there are. To do this, we first need a lemma:

2.7 Lemma *If $f(x)$ is an irreducible polynomial of degree n over \mathbf{Z}_p , then $f(x)|x^{p^n} - x$.*

PROOF. By adjoining a root of f if necessary, we can assume that $f(x)$ has a root α in an extension field K of \mathbf{Z}_p . We have already seen that $\mathbf{Z}_p[\alpha]$ is then a finite field with p^n elements, and that every element of this field (including α) is a root of the polynomial $x^{p^n} - x$. But then Corollary 2.2 shows that $f(x)|x^{p^n} - x$. \square

2.8 Theorem *For each prime p and each positive integer n , there is (up to isomorphism) exactly one field with p^n elements.*

PROOF. We have already seen that there is *at least* one such field.

Let f be any irreducible polynomial over \mathbf{Z}_p of degree n . We know that at least one such polynomial exists.

Now say K is a field with p^n elements. First, we know that the characteristic of K must be p . For the number of elements in K is a power of the characteristic.

We know that the multiplicative group K^* of this field is cyclic, so every element of this field is a root of the polynomial $x^{p^n} - x$. Now let us factor $x^{p^n} - x$ into irreducible polynomials over \mathbf{Z}_p : we have

$$x^{p^n} - x = g_1(x)g_2(x)\cdots g_r(x)$$

We know that

$$\sum_{i=1}^r \deg(g_i) = p^n$$

We also know that one of these polynomials must be $f(x)$; say $g_1(x)$ is $f(x)$. Now f must have a root α in K , for if not, the number of roots of $x^{p^n} - x$ would only be at most

$$\sum_{i=2}^r \deg g_i < p^n$$

and this is a contradiction.

So α is in K , and therefore $\mathbf{Z}_p[\alpha]$ (which is just the set of all polynomials in α with coefficients in K) is also contained in K . Since both K and $\mathbf{Z}_p[\alpha]$ have p^n elements, we must have $\mathbf{Z}_p[\alpha] = K$. Thus, K is derived from an irreducible polynomial over \mathbf{Z}_p of degree n , and any such polynomial gives rise to K , so they all generate isomorphic fields. \square

Now we know that any finite field has to have p^n elements, where p is the characteristic of the field and n is its dimension as a vector space over its prime field. So we have just seen that there is (up to isomorphism) exactly one such field for each p and n , and it can be constructed from an irreducible polynomial of degree n over \mathbf{Z}_p .

The standard notation for a field with q elements is F_q , so we have just constructed F_{p^n} , and we have shown that it is unique.

2.4 Primitive polynomials

A *primitive polynomial* over a field K is a polynomial p with coefficients in K (i.e., p is in $K[x]$) that is irreducible over K and is such that if α is a root of p in an extension field of K , then α is a primitive element of $K[\alpha]$.

That is, p is a primitive polynomial over K if and only if the powers of α are all the non-zero elements of $K[\alpha]$.

This definition makes sense, since as we saw previously, all the roots of an irreducible polynomial have the same multiplicative order.

Not every irreducible polynomial is primitive. On the other hand, the construction at the end of the last section shows that there are primitive polynomials of every order over \mathbf{Z}_p . Actually, the same proof shows that there are primitive polynomials of every order over every finite field.

Here is a simple example: Let us consider polynomials over \mathbf{Z}_p . Every polynomial of degree 1 (i.e., of the form $p(x) = ax + b$ with a and b in \mathbf{Z}_p and $a \neq 0$) is irreducible. That's always true for polynomials over any field—non-trivial first degree polynomials can't be factored. On the other hand, the polynomial $p(x) = x - a$ is a primitive polynomial if and only if a is a primitive element of \mathbf{Z}_p . So for instance, $x - 3$ and $x - 5$ are primitive polynomials over \mathbf{Z}_7 , while $x - a$ (for a different from 3 and 5) are irreducible but not primitive over \mathbf{Z}_7 . This is really a degenerate case—the root 3 of $x - 3$ is in the base field \mathbf{Z}_7 itself, not in an extension field. The actual definition above glossed over this point, although it wasn't actually incorrect.

However, a more convincing example can be arrived at by looking at polynomials

of degree 4 over \mathbf{Z}_2 . Let us consider the two polynomials

$$\begin{aligned} f(x) &= x^4 + x + 1 \\ g(x) &= x^4 + x^3 + x^2 + x + 1 \end{aligned}$$

They are both irreducible polynomials. We won't bother to prove this, but of course in principle one could verify this by just trying all possible factorizations; there are only a (small) finite number of them. Therefore they each generate a 4-dimensional field over \mathbf{Z}_2 , and as we have stated, these two fields must be the same except for renaming—each is (isomorphic to) $F_{2^4} = F_{16}$. So we can assume that they are actually the same field.

Now it turns out that f is primitive over \mathbf{Z}_2 , while g is not. It's easy to see this. Suppose that α is a root of f and β is a root of g . Then just by using the polynomial relations $f(\alpha) = 0$ and $g(\beta) = 0$, we can construct the following tables of powers of α and β :

n	α^n	β^n
0	1	1
1	α	β
2	α^2	β^2
3	α^3	β^3
4	$\alpha + 1$	$\beta^3 + \beta^2 + \beta + 1$
5	$\alpha^2 + \alpha$	1
6	$\alpha^3 + \alpha^2$	
7	$\alpha^3 + \alpha + 1$	
8	$\alpha^2 + 1$	
9	$\alpha^3 + \alpha$	
10	$\alpha^2 + \alpha + 1$	
11	$\alpha^3 + \alpha^2 + \alpha$	
12	$\alpha^3 + \alpha^2 + \alpha + 1$	
13	$\alpha^3 + \alpha^2 + 1$	
14	$\alpha^3 + 1$	
15	1	

We can see that α has order $15 = 2^4 - 1$, while β has order 5. Since the non-zero elements of F_{16} are just the powers of α , it must be the case that β is a power of α . (In particular, this is why the order of β has to be a divisor of the order of α .) Any power of α that has order 5 will do. For instance, α^3 clearly has order

5. And, using the above table of powers of α , we can see that $g(\alpha^3) = 0$:

$$\begin{aligned}
 g(\alpha^3) &= (\alpha^3)^4 + (\alpha^3)^3 + (\alpha^3)^2 + \alpha^3 + 1 \\
 &= \alpha^{12} + \alpha^9 + \alpha^6 + \alpha^3 + 1 \\
 &= \alpha^3 + \alpha^2 + \alpha + 1 && (= \beta^4) \\
 &\quad + \alpha^3 && + \alpha && (= \beta^3) \\
 &\quad + \alpha^3 + \alpha^2 && && (= \beta^2) \\
 &\quad + \alpha^3 && && (= \beta) \\
 &\quad + && && 1 && (= 1) \\
 &= 0
 \end{aligned}$$

Note what is going on here: setting $\beta = \alpha^3$ creates an isomorphism between $\mathbf{Z}_2[\beta]$ and $\mathbf{Z}_2[\alpha]$. There are other isomorphisms as well—for instance, $\beta = \alpha^6$ would also work. But once we have picked such an isomorphism, we don't change it.

It's important to bear in mind that α and β *both* generate F_{16} algebraically. If as above, we take $\beta = \alpha^3$, then we can see as in the above computation that

$$\alpha = \beta^4 + \beta^2 + 1$$

and so, substituting for β^4 and simplifying,

$$\alpha = \beta^3 + \beta$$

So α , and therefore any power of α , is a polynomial in β . Thus even though β is not a primitive element of F_{16} , it certainly generates it—it's just that we need more polynomials in β than just its powers to get all the elements of F_{16} .

Chapter 3

Some Random Number Generators

In this chapter we give a brief survey of some random number generators that make use of the theory of finite fields.

3.1 Linear pseudo-congruential generators

A linear pseudo-congruential generator is defined by choosing

- a prime number p ,
- integers $a \geq 1$ and $b \geq 0$,
- an initial value (or “seed”) s_0 in the range $0 \leq s_0 < p$,

and then recursively computing

$$s_n = as_{n-1} + b \pmod{p}$$

This notation, by the way, always means that the right-hand side is replaced by the number in the range 1 to $p - 1$ to which it is congruent (mod p). Such a sequence is always ultimately periodic. In fact, if $s_i = s_j$ for some i and j , then $s_{i+1} = s_{j+1}$, and so on, so that the sequence repeats from that point on. This also shows that the period length cannot be greater than p . The period length can of course be p —just set $a = 1$ and $b = 1$. This is not a useful source of pseudo-random numbers, however. What makes a pseudo-random number generator useful is an involved topic. Here we shall just confine ourselves to finding generators that have a long period length.

Now it might seem that picking b cleverly would add to the complexity, or at least the period length of the sequence, but this is really not the case. In fact, we can clearly get a period length of $p - 1$ by letting $b = 0$, choosing s_0 in the range $1 \leq s_0 < p$ (i.e., making sure that $s_0 \neq 0$), and choosing for a any primitive root (mod p). The numbers s_n will then just be the powers of a (mod p), which are all the numbers from 1 to $p - 1$ (in a permuted order, of course).

3.2 Higher order linear recursive generators

This is similar, except that we use a higher-order linear recursion. To form a recursive generator of degree r , we start by picking r initial seeds s_0, s_1, \dots, s_{r-1} and constants c_0, c_1, \dots, c_{r-1} . The seeds are in the range 0 to $p - 1$, but are not all 0, and the constants are in the range 0 to $p - 1$ (and of course they are also not all 0). Then for $n \geq r$, we compute recursively

$$s_n = c_{r-1}s_{n-1} + c_{r-2}s_{n-2} + \dots + c_0s_{n-r} \pmod{p}$$

Thus, each successive value of s_n is a linear combination (mod p) of the previous r values, where the coefficients are constant. This enables us, as we will see, to get a much larger period length.

To avoid trivial complications, we will assume that $c_0 \neq 0$. (This is really no restriction; if it is zero, we just have a recursive generator of degree $r - 1$ instead of r .)

We define a polynomial in $\mathbf{Z}_p[x]$ of degree r corresponding to such a higher order generator by

$$f(x) = x^r - c_{r-1}x^{r-1} - c_{r-2}x^{r-2} - \dots - c_0$$

The theory of this recursion turns out to be simplest in the case that f is irreducible over \mathbf{Z}_p ; we will assume that this is the case.

We have to clarify exactly what we mean by the term “period length” for our sequence $\{s_j\}$. But once we do, it will turn out that the period length can be no more than $p^r - 1$, and it will have this value precisely when f is a primitive polynomial (mod p). The rest of this section is devoted to proving this result.

1. We define a sequence of vectors for $n \geq 0$:

$$v_n = \langle s_{n+r-1}, s_{n+r-2}, \dots, s_n \rangle$$

That is, v_n is just the vector in \mathbf{Z}_p^r whose components are r consecutive terms in the sequence $\{s_i\}$. Each v_n moves down this sequence 1 element. Choosing the initial set of seeds s_0, s_1, \dots, s_{r-1} is just equivalent to choosing an initial vector $v_0 = \langle s_{r-1}, s_{r-2}, \dots, s_0 \rangle$.

Now let us say precisely what we mean by the period length of the sequence $\{s_i\}$. It is perfectly possible to have $s_3 = s_8$, for instance, without this forcing $s_4 = s_9$. This is because each element s_j of the sequence depends, not just on the previous element of the sequence, but on the previous r elements. On the other hand, if $v_3 = v_8$, for instance, then we know by the recursive formula that v_4 must also equal v_9 , and in general $v_i = v_{i+5}$ for all $i \geq 3$. So when we talk about the period length of the sequence generated by our recursive formula, what we really mean is the period length of the sequence of vectors $\{v_i\}$.

Clearly, v_{n+1} can be computed from v_n ; in fact,

$$v_{n+1} = \left\langle \sum_{i=1}^r c_{r-i} s_{n+r-i}, s_{n+r-1}, s_{n+r-2}, \dots, s_{n+1} \right\rangle$$

where all the components have been shifted to the right one position, the right-most one discarded, and a new component (computed by the recursive formula) placed in the left-most position.

This is a linear transformation on \mathbf{Z}_p^r . Let us denote it by the symbol R (for “recursion”). R has the matrix representation

$$R = \begin{pmatrix} c_{r-1} & c_{r-2} & c_{r-3} & \dots & c_1 & c_0 \\ 1 & 0 & 0 & \dots & 0 & 0 \\ 0 & 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 1 & 0 \end{pmatrix}$$

(Notice that all the elements on the main diagonal are 0 except for the upper left-hand element.)

Now R is an invertible matrix—its determinant¹ is $(-1)^{r-1}c_0 \neq 0$. Hence the inverse of R —call it R^{-1} exists. In fact, then, all powers R^n of R exist, where n can be any integer (positive, negative, or 0). This set of powers of R forms a group. It is a subset of the set of $r \times r$ matrices with elements in \mathbf{Z}_p . There are only a finite number of such matrices. Therefore, the same argument as at

¹I know, I didn't talk about determinants. Trust me. Either that, or solve for v_n in terms of v_{n+1} to get an explicit construction for R^{-1} :

$$R^{-1} = \begin{pmatrix} 0 & 1 & 0 & \dots & 0 & 0 \\ 0 & 0 & 1 & \dots & 0 & 0 \\ 0 & 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 0 & 1 \\ \frac{1}{c_0} & -\frac{c_{r-1}}{c_0} & -\frac{c_{r-2}}{c_0} & \dots & -\frac{c_2}{c_0} & -\frac{c_1}{c_0} \end{pmatrix}$$

the beginning of Section 1.3.4 shows that $\{R^i\}$ is a finite cyclic group. Say its order is m . (Usually we just say that the order of R is m ; this just means that R has order m in the cyclic group of powers of R .)

Now we have $v_{n+1} = Rv_n$, and in general, $v_n = R^n v_0$. Thus, if R has order m , then the period of the sequence of vectors $\{v_j\}$ is $\leq m$. Further, since R is invertible, as long as we start from a vector v_0 that is not 0 (i.e., as long as not all of s_0, s_1, \dots, s_{r-1} are 0), none of the vectors v_j will be the 0 vector. Thus, the period length of the sequence $\{v_j\}$ can be at most $p^r - 1$.

We will show that the period length of the sequence $\{v_j\}$ is $p^r - 1$ if and only if the order of R is $p^r - 1$.

Suppose first that the period length of the sequence $\{v_j\}$ is $p^r - 1$. Then the set of vectors v_0, v_1, \dots is all of $(\mathbf{Z}_p^r)^*$, since there are just $p^r - 1$ elements of $(\mathbf{Z}_p^r)^*$. Thus, if v is any non-zero vector in \mathbf{Z}_p^r , we have $R^{p^r-1}v = v$. Since this holds for all non-zero vectors, we must have $R^{p^r-1} = I$; that is, the order of R is $\leq p^r - 1$. On the other hand, we just saw that the order of R is \geq the period length of the sequence $\{v_j\}$, so in fact the order of R is $p^r - 1$.

Now we go in the other direction. Suppose that the order of R is $p^r - 1$. We must show that the period length of the sequence is also $p^r - 1$. Let us pick a new initial vector

$$w_0 = \langle 1, 0, 0, \dots, 0 \rangle$$

We can see from the recursive formula that

$$\begin{aligned} w_0 &= \langle 1, 0, 0, \dots, 0 \rangle \\ w_1 &= \langle *, 1, 0, \dots, 0 \rangle \\ w_2 &= \langle *, *, 1, \dots, 0 \rangle \\ &\vdots \\ w_{r-1} &= \langle *, *, *, \dots, 1 \rangle \end{aligned}$$

where “*” denotes some expression, which may or may not be zero—we don’t care.

Now the point is this: these vectors w_0, w_1, \dots, w_{r-1} span all of \mathbf{Z}_p^r . This is because w_0 is the “unit vector in the x direction”. We then get the “unit vector in the y” direction by subtracting some multiple of w_0 from w_1 . We get the “unit vector in the z direction” by subtracting multiples of w_0 and w_1 from w_2 ; and so on. In this way we get r vectors that clearly span the space, and so our original vectors do also, since the “unit vectors” can be expressed in terms of them.

Now suppose the period length of the sequence $\{v_i\}$ was less than $p^r - 1$. Then also the period length of the sequence $\{w_i\}$ starting with w_0 would be less than

$p^r - 1$. (Otherwise, the sequence starting with w_0 would be all of $(\mathbf{Z}_p^r)^*$, and would include the original sequence; and that in turn would mean that the original sequence would have period length $p^r - 1$, which would be a contradiction. So say the period length of the sequence starting with w_0 is $m < p^r - 1$. We know that $R^m w_0 = w_0$, and in fact we have

$$R^m w_i = R^m R^i w_0 = R^{m+i} w_0 = R^i R^m w_0 = R^i w_0 = w_i$$

Thus $R^m w_i = w_i$ on the spanning set of vectors w_0, w_1, \dots, w_{r-1} . Hence $R^m = I$ on all of \mathbf{Z}_p^r . This is a contradiction, since we have assumed that the order of R was $p^r - 1$.

Thus we have in fact showed that the order of R is $p^r - 1$ if and only if the period length of the sequence $\{v_i\}$ is $p^r - 1$.

2. Now since f is irreducible² of degree r , the field F_{p^r} is just the field $\mathbf{Z}_p[\alpha]$ obtained by adjoining a root α of f to \mathbf{Z}_p . Let us look at the linear operator M defined on F_{p^r} by “multiplication by α ”. That is, we define

$$Mv = \alpha v$$

for all v in F_{p^r} . Since F_{p^r} is as a vector space equivalent to \mathbf{Z}_p^r , M has a representation as a matrix, and this representation is just

$$M = \begin{pmatrix} c_{r-1} & 1 & 0 & \dots & 0 & 0 \\ c_{r-2} & 0 & 1 & \dots & 0 & 0 \\ c_{r-3} & 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ c_1 & 0 & 0 & \dots & 0 & 1 \\ c_0 & 0 & 0 & \dots & 0 & 0 \end{pmatrix}$$

We see that M is just the adjoint of R . (This is the key to the proof.) Since $M^m = I$ if and only if $R^m = (M^*)^m = (M^m)^* = I^* = I$, we see that the order of M is just the order of R .

Finally, the order of M is just $p^r - 1$ if and only if the root α of f is a primitive root, i.e., if and only if the polynomial f is a primitive polynomial of degree r over \mathbf{Z}_p .

3. Putting these results all together, we have shown that provided f is an irreducible polynomial over \mathbf{Z}_p ,

$$\begin{aligned} f \text{ is a primitive polynomial over } \mathbf{Z}_p &\text{ iff } M \text{ has order } p^r - 1 \\ &\text{ iff } R \text{ has order } p^r - 1 \\ &\text{ iff the period length of } \{v_i\} \text{ is } p^r - 1 \end{aligned}$$

and the proof is complete.

²This is the only place where the fact that f is irreducible is needed.

3.3 Inversive congruential generators

An inversive congruential generator is a non-linear recursive generator having the form

(3.1)

$$s_{n+1} = \begin{cases} b & \text{if } s_n = 0 \\ as_n^{-1} + b \pmod{p} & \text{if } s_n \neq 0 \end{cases}$$

where a and b are constants, $a \not\equiv 0 \pmod{p}$, and s_n^{-1} is the multiplicative inverse of s_n in \mathbf{Z}_p . Of course we have to initialize the sequence; we will do so by setting

$$s_0 = b$$

Note that we have to treat the case when $s_n = 0$ separately, since of course 0 does not have a multiplicative inverse in \mathbf{Z}_p . Setting s_{n+1} to b in this case is the only reasonable choice, since b cannot be the value of $as_n + b$ in any other case.

Since each term depends only on the previous one, we see that the maximum possible period length for such a sequence is p . Flahive and Niederreiter, extending earlier work of Eichenauer and Lehn, have given a complete characterization of generators of this form that have period p . Here is how:

First of all, let us associate the polynomial

$$f(x) = x^2 - bx - a$$

with this recursion. (That is, the coefficients a and b in this quadratic polynomial are determined from the recursive definition.) There are two possibilities:

Case I: f is reducible. In this case, f factors into 2 linear factors in $\mathbf{Z}_p[x]$:

$$f(x) = (x - r)(x - s)$$

We know that not both of r and s are 0, since otherwise both a and b would be 0, and we have assumed that $a \neq 0$. So say $r \neq 0$. Since $f(r) = 0$, we have

$$r^2 - br - a = 0$$

in \mathbf{Z}_p . Letting r^{-1} denote the multiplicative inverse of r in \mathbf{Z}_p and multiplying through by r^{-1} , we get

$$r = ar^{-1} + b$$

But this means that the recursion cannot have order p , since the only way it could have order p is to visit each element of \mathbf{Z}_p over and over again,

but what we have just shown is that once the sequence hits the value r , it has to stay there.

Therefore, the sequence can only have maximal period p if f is irreducible. (Not every irreducible f yields a maximal period, however, as we will see below.)

Case II: f is irreducible. In this case, f does not factor over \mathbf{Z}_p , but we know that it does factor over F_{p^2} . Say it factors as $f(x) = (x - \alpha)(x - \beta)$ where α and β are in F_{p^2} . We can show that in this case, α and β must be distinct. For suppose they were the same. Then we would have

$$f(x) = \begin{cases} x^2 - bx - a \\ (x - \alpha)^2 = x^2 - 2\alpha x + \alpha^2 \end{cases}$$

so we would have to have $b = 2\alpha$. There are two cases to consider:

The characteristic p is not 2. In this case, 2 has a multiplicative inverse in \mathbf{Z}_p , and we get $\alpha = b/2$. That is, α is already in \mathbf{Z}_p , which means that f actually is reducible over \mathbf{Z}_p —but we have assumed this is not true.

The characteristic p is 2. In this case, since $b = 2\alpha$, b must be 0. Hence $f(x) = x^2 - a$. There are only two possibilities:

- $f(x) = x^2$. This is of course reducible.
- $f(x) = x^2 + 1 = (x + 1)^2$. And this is also reducible.

So we see that if f is irreducible, its roots in F_{p^2} must be distinct.

We have seen so far that an inversive congruential generator can only have maximal period if its associated polynomial f is irreducible, and in that case, its two roots in F_{p^2} must be distinct.

The next thing we do is linearize the recursion by making the substitution
(3.2)

$$s_n = \frac{t_{n+1}}{t_n}$$

Thus, so long as $t_n \neq 0$,

$$s_n = \begin{cases} \frac{t_{n+1}}{t_n} \\ a \frac{t_{n-1}}{t_n} + b \end{cases}$$

and so we have

(3.3)

$$t_{n+1} = bt_n + at_{n-1}$$

This substitution may seem like just a clever unmotivated trick. Flahive and Niederreiter attempt to motivate it by pointing out that the original non-linear recursion (3.1) is essentially that encountered in the theory of continued fractions. For our purposes, we don't need to know that.

Now (3.3) is a second-order linear recursion. We have already seen that there is a nice theory of such recursions in Section 1.6.4 (page 33), where we used it to find a formula for the n^{th} term of the Fibonacci sequence. Exactly the same reasoning applies here, even though now we are considering equations and vector spaces over \mathbf{Z}_p , instead of over \mathbf{R} . Just as in Section 1.6.4, we know that there are two special solutions of this recursion:

- $1, \alpha, \alpha^2, \alpha^3, \dots$
- $1, \beta, \beta^2, \beta^3, \dots$

where α and β are the two distinct roots of $f(x)$ in F_{p^2} . Further, any solution of the recursive equation 3.3 can be written as a linear combination (over \mathbf{Z}_p) of these two solutions. So the special solution s_n can be written as

$$t_n = c_1\alpha^n + c_2\beta^n$$

where c_1 and c_2 have to be determined. They are determined by using the facts that $s_0 = b$. This forces $t_0 = 1$ and $t_1 = b$.³ So we get

$$\begin{aligned} c_1 + c_2 &= 1 \\ c_1\alpha + c_2\beta &= b \end{aligned}$$

Solving these two equations for c_1 and c_2 yields

$$\begin{aligned} c_1 &= \frac{\alpha}{\alpha - \beta} \\ c_2 &= -\frac{\beta}{\alpha - \beta} \end{aligned}$$

and so
(3.4)

$$t_n = \frac{\alpha^{n+1} - \beta^{n+1}}{\alpha - \beta}$$

Consecutive values of s_n can be computed from this formula and 3.2 until the point at which t_n first becomes 0.

The crux of the matter is now contained in the following lemma, which we have already done most of the work of proving:

³Actually, any non-zero value for t_0 would do, and then t_1 would just be b times that value. So there is no harm in taking $t_0 = 1$; any other value would just multiply all the terms t_n by that value and leave the values of the terms s_n unchanged.

3.1 Lemma *Let $f(x) = x^2 - bx - a$ be irreducible over \mathbf{Z}_p , and let α and β be its roots in F_{p^2} . Let N be the order of the element β/α in the multiplicative group $F_{p^2}^*$. Then the period length of the sequence s_n generated by (3.1) with $x_0 = b$ is $N - 1$.*

PROOF. For all n such that $1 \leq n < N$, $\beta^n \neq \alpha^n$, while $\beta^N = \alpha^N$. Equation (3.4) then shows that $t_n \neq 0$ for all n such that $0 \leq n < N-1$, and $t_{N-1} = 0$. Equation (3.2) then shows that $s_n \neq 0$ for all n such that $0 \leq n < N-2$ and $s_{N-2} = 0$. Equation (3.1) then shows that $s_n \neq b$ for all n such that $1 \leq n < N-1$, and $s_{N-1} = b$, and we are done, since we know that also $s_0 = b$. \square

3.2 Theorem *If the polynomial f corresponding to an inversive congruential generator has roots α and β in F_{p^2} , the period length of the generator is maximal if and only if the order of β/α in the multiplicative group $F_{p^2}^*$ is $p + 1$.*

PROOF. 1. If $\{s_n\}$ has a maximal period length, then we have already seen that f is irreducible, and the lemma then shows that the order of α/β is $p + 1$.

2. If the order of β/α is $p + 1$, then f must be irreducible, since otherwise both β and α (and hence also β/α) would be in \mathbf{Z}_p and the order of β/α would divide $p - 1$. Hence the lemma applies again. \square

Flahive and Niederreiter say that a polynomial $f(x) = x^2 - bx - a$ is an *inversive maximal period polynomial* (or an *IMP polynomial*) if the period length of the associated sequence $\{s_n\}$ is p . The theorem just proved thus characterizes IMP polynomials. In particular, it turns out that all primitive monic quadratic polynomials are IMP polynomials:

3.3 Corollary *If f is a primitive polynomial, then f is an IMP polynomial.*

PROOF. If α is a root of f as above, then every element of $F_{p^2}^*$ is a power of α . So in particular, $\beta = \alpha^t$ for some t such that $1 < t < p^2 - 1$. (t can't be 1 because we know that $\beta \neq \alpha$, and t can't be $p^2 - 1$ because we know that $\beta \neq 1$.) In fact, we can show that t must be p . The reason is as follows: we have seen that F_{p^2} is isomorphic to the fields $\mathbf{Z}_p[\alpha]$ and $\mathbf{Z}_p[\beta]$. Now in $\mathbf{Z}_p[\alpha]$, we know that α^t is also a root of f . This means that in $\mathbf{Z}_p[\beta]$ (and therefore also in F_{p^2}), β^t is also a root of f . But if $t \neq p$, then the elements α , α^t (which is β), and α^{t^2} (which is β^t) are all distinct, and so f would have to have three roots, which we know is impossible. Therefore $t = p$ and $\beta = \alpha^p$.

So we have $\beta/\alpha = \alpha^{p-1}$. If then $(\beta/\alpha)^n = 1$, we have $\alpha^{n(p-1)} = 1$. Since α is a primitive root, this is true precisely when

$$p^2 - 1 | n(p - 1)$$

Since $p^2 - 1 = (p + 1)(p - 1)$, this is the same as saying $p + 1 | n$. Thus the order of β/α is $p + 1$, and so by the theorem, f is an IMP polynomial. \square

Bibliography

- [1] Emil Artin. *Galois Theory*. University of Notre Dame, Notre Dame, Indiana, 1942. [QA1.N87 no. 2].
- [2] Michael Artin. *Algebra*. Prentice Hall, Engelwood Cliffs, New Jersey, 1991. [QA154.2.A77].
- [3] Mary Flahive and Harald Niederreiter. On Inversive Congruential Generators for Pseudorandom Numbers. In Gary L. Mullen and Peter Jau-Shyong Shiue, editors, *Finite Fields, Coding Theory, and Advances in Communications and Computing*, pages 75–80. M. Dekker, New York, 1992. [QA247.3.F56].
- [4] G. H. (Godfrey Harold) Hardy and E. M. Wright. *An Introduction to the Theory of Numbers*. Oxford University Press, New York, 1979. [QA241.H28].
- [5] Rudolf Lidl and Harald Niederreiter. *Finite Fields*. Cambridge University Press, Cambridge (England); New York, 2nd edition, 1997. [QA247.3.L53].
- [6] Neal Zierler. Linear Recurring Sequences. *Journal of the Society for Industrial and Applied Mathematics*, 7(1):31–48, March 1959.