

MA251: Algebra I – Advanced Linear Algebra

Christian Böhning

Based on notes of Adam Thomas, Derek Holt and David Loeffler

Term 1, 2022-23

Contents

0	Blurb	2
1	Review of some MA106 material	3
1.1	Fields	3
1.2	Vector spaces	4
1.3	Linear maps	5
1.4	The matrix of a linear map with respect to a choice of (ordered) bases	6
1.5	Change of basis	7
2	The Jordan Canonical Form	11
2.1	Introduction	11
2.2	Eigenvalues and eigenvectors	11
2.3	The minimal polynomial	12
2.4	The Cayley–Hamilton theorem	14
2.5	Calculating the minimal polynomial	15
2.6	Jordan chains and Jordan blocks	18
2.7	Jordan bases and the Jordan canonical form	20
2.8	The JCF when $n=2$ and 3	22
2.9	The general case	25
2.10	Examples	26
2.11	Proof of Theorem 2.7.3	28
2.12	An algorithm to compute the Jordan canonical form in general (brute force)	30
2.13	Grand finale	31
3	Functions of matrices	33
3.1	Powers of matrices	33
3.2	Applications to difference equations	35
3.3	Motivation: Systems of Differential Equations	37
3.4	Definition of a function of a matrix	38

4	Bilinear Maps and Quadratic Forms	42
4.1	Bilinear maps: definitions	42
4.2	Bilinear maps: change of basis	43
4.3	Quadratic forms	45
4.4	Nice bases for quadratic forms	47
4.5	Euclidean spaces, orthonormal bases and the Gram–Schmidt process	53
4.6	Orthogonal transformations	55
4.7	Nice orthonormal bases	58
4.8	Applications of quadratic forms to geometry	63
4.8.1	Reduction of the general second degree equation	63
4.8.2	The case $n = 2$	65
4.8.3	The case $n = 3$	65
4.9	Singular value decomposition	72
4.10	The complex story	74
4.10.1	Sesquilinear forms	74
4.10.2	Operators on Hilbert spaces	76
5	Finitely Generated Abelian Groups	78
5.1	Definitions	78
5.2	Subgroups, cosets and quotient groups	80
5.3	Homomorphisms and the first isomorphism theorem	83
5.4	Free abelian groups	85
5.5	Unimodular elementary row and column operations and the unimodular Smith normal form for integer matrices	87
5.6	Subgroups of free abelian groups	92
5.7	General finitely generated abelian groups	94
5.8	Finite abelian groups	96

0 Blurb

As its title suggests, this module is a continuation of last year’s MA106 Linear Algebra module; we’ll be studying vector spaces, linear maps, and their properties in a bit more detail. Later in the module we’ll think a bit about matrices whose entries lie not in a field but in the integers \mathbb{Z} , and we’ll see what our methods have to tell us in that case.

1 Review of some MA106 material

In this section, we'll recall some ideas from the first year MA106 Linear Algebra module. This will just be a brief reminder; for detailed statements and proofs, go back to your MA106 notes.

1.1 Fields

Recall that a *field* is a number system where we know how to do all of the basic arithmetic operations: we can add, subtract, multiply and divide (as long as we're not trying to divide by zero).

Definition 1.1.1. A field is a non-empty set K together with two operations (maps from $K \times K$ to K) addition, denoted by $+$, and multiplication, denoted by \cdot (or just juxtaposition), satisfying the following axioms:

1. $a + b = b + a$ for all $a, b \in K$;
2. there exists an element $0 \in K$ such that $a + 0 = a$ for all $a \in K$;
3. $(a + b) + c = a + (b + c)$ for all $a, b, c \in K$;
4. there exists an element $-a \in K$ such that $a + (-a) = 0$ for all $a \in K$;
5. $a \cdot b = b \cdot a$;
6. there exists an element $1 \in K, 1 \neq 0$, such that $1 \cdot a = a$ for all $a \in K$;
7. $(a \cdot b) \cdot c = a \cdot (b \cdot c)$ for all $a, b, c \in K$;
8. there exists an element $a^{-1} \in K$ such that $a \cdot a^{-1} = 1$ for all $0 \neq a \in K$;
9. $a \cdot (b + c) = (a \cdot b) + (a \cdot c)$ for all $a, b, c \in K$.

Examples.

- A non-example is \mathbb{Z} , the integers. Here we can add, subtract, and multiply, but we can't always divide without jumping out of \mathbb{Z} into some bigger world. That is to say that Axiom 8 would fail: there are no multiplicative inverses of any integer apart from 1 and -1 .
- The real numbers \mathbb{R} and the complex numbers \mathbb{C} are fields, and these are perhaps the most familiar ones.
- The rational numbers \mathbb{Q} are also a field.
- A more subtle example: if p is a prime number, the integers mod p are a field, written as $\mathbb{Z}/p\mathbb{Z}$ or \mathbb{F}_p .

There are lots of fields out there, and the reason we take the axiomatic approach is that we know that everything we prove will be applicable to any field we like, as long as we've only used the field axioms in our proofs (rather than any specific properties of the fields we happen to most

1 Review of some MA106 material

like). We don't have to know all the fields in existence and check that our proofs are valid for each one separately.

1.2 Vector spaces

Let K be a field¹. A *vector space* over K is a non-empty set V together with two extra pieces of structure. Firstly, it has to have a notion of *addition*: we need to know what $v + w$ means if v and w are in V . Secondly, it has to have a notion of *scalar multiplication*: we need to know what λv means if v is in V and λ is in K . These have to satisfy some axioms, for which I'm going to refer you again to your MA106 notes.

Definition 1.2.1. A vector space V over a field K is a set V with two operations. The first is addition, a map from $V \times V$ to V satisfying Axioms 1 to 4 in the definition of a field. The second operation is scalar multiplication, a map from $K \times V$ to V denoted by juxtaposition or \cdot , satisfying the following axioms:

1. $\alpha(u + v) = \alpha u + \alpha v$ for all $u, v \in V, \alpha \in K$;
2. $(\alpha + \beta)v = \alpha v + \beta v$ for all $v \in V, \alpha, \beta \in K$;
3. $(\alpha \cdot \beta)v = \alpha(\beta v)$ for all $v \in V, \alpha, \beta \in K$;
4. $1 \cdot v = v$ for all $v \in V$.

A *basis* of a vector space is a subset $B \subset V$ such that every $v \in V$ can be written *uniquely* as a finite linear combination of elements of B ,

$$v = \lambda_1 b_1 + \cdots + \lambda_n b_n,$$

for some $n \in \mathbb{N}$ and some $\lambda_1, \dots, \lambda_n \in K$. So for each $v \in V$, we can do this in one and only one way. Another way of saying this is that B is a linearly independent set which spans V , which is the definition you had in MA106. We say V is *finite-dimensional* if there is a finite basis of V . You saw last year that if V has one basis which is finite, then every basis of V is finite, and they all have the same cardinality; and we define the *dimension* of V to be this number which is the number of elements in any basis of V .

Examples. Let $K = \mathbb{R}$.

- The space of polynomials in x with coefficients in \mathbb{R} is certainly a vector space over \mathbb{R} ; but it's not finite-dimensional (rather obviously).
- For any $d \in \mathbb{N}$, the set \mathbb{R}^d of column vectors with d real entries is a vector space over \mathbb{R} (which, not surprisingly, has dimension d).
- The set

$$\left\{ \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \in \mathbb{R}^3 : x_1 + x_2 + x_3 = 0 \right\}$$

¹It's conventional to use K as the letter to denote a field; the K stands for the German word "Körper".

1 Review of some MA106 material

is a vector space over \mathbb{R} if we define vector addition and scalar multiplication component-wise as usual.

The third example above is an interesting one because there's no "natural choice" of basis. It certainly has bases, e.g. the set

$$\left\{ \begin{pmatrix} 1 \\ -2 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \right\},$$

but there's no reason why that's better than any other one. This is one of the reasons why we need to worry about the choice of basis – if you want to tell someone else all the wonderful things you've found out about this vector space, you might get into a total muddle if you insisted on using one particular basis and they preferred another different one.

The following lemma (which will be required in the proof of one of our main theorems) is straightforward from the material in MA106 - the proof is left as an exercise to check you are comfortable with such material.

Lemma 1.2.2. *Suppose that U is an m -dimensional subspace of an n -dimensional vector space V and $\mathbf{w}_1, \dots, \mathbf{w}_{n-m}$ extend a basis of U to a basis of V . Then the equation*

$$\alpha_1 \mathbf{w}_1 + \dots + \alpha_{n-m} \mathbf{w}_{n-m} + \mathbf{u} = \mathbf{0}, \quad \text{where } \mathbf{u} \in U, \quad (1)$$

only has the solution $\alpha_i = 0$ for all $1 \leq i \leq n - m$ and $\mathbf{u} = \mathbf{0}$.

1.3 Linear maps

If V and W are vector spaces (over the same field K), then a *linear map* from V to W is a map $T : V \rightarrow W$ which "respects the vector space structures". That is, we know two things that we can do with vectors in a vector space – add them, and multiply them by scalars; and a linear map is a map where adding or scalar-multiplying on the V side, then applying the map T , is the same as applying the map T , then adding or multiplying on the W side. Formally, for T to be a linear map means that we must have

$$T(v_1 + v_2) = T(v_1) + T(v_2) \quad \forall v_1, v_2 \in V$$

and

$$T(\lambda v_1) = \lambda T(v_1) \quad \forall \lambda \in K, v_1 \in V.$$

Example 1. Let V and W be vector spaces over K . Then $T : V \rightarrow W$ defined by $T(v) = 0_W = \mathbf{0}$ for all $v \in V$ is a linear map, called the zero linear map. Furthermore, we have $S : V \rightarrow V$ defined by $S(v) = v$ for all $v \in V$ is a linear map, called the identity linear map.

Example 2. Let $V = \mathbb{R}^3$ and $W = \mathbb{R}^2$. Then the following maps $T : V \rightarrow W$ are linear.

$$1. T \left(\begin{pmatrix} a \\ b \\ c \end{pmatrix} \right) = \begin{pmatrix} a \\ b \end{pmatrix};$$

1 Review of some MA106 material

$$2. T \left(\begin{pmatrix} a \\ b \\ c \end{pmatrix} \right) = \begin{pmatrix} b \\ 0 \end{pmatrix};$$

$$3. T \left(\begin{pmatrix} a \\ b \\ c \end{pmatrix} \right) = \begin{pmatrix} a+b \\ b+c \end{pmatrix}.$$

Whereas, you should check that $T \left(\begin{pmatrix} a \\ b \\ c \end{pmatrix} \right) = \begin{pmatrix} a^2 \\ b \end{pmatrix}$ is NOT a linear map.

1.4 The matrix of a linear map with respect to a choice of (ordered) bases

Let V and W be vector spaces over a field K . Let $T : V \rightarrow W$ be a linear map, where $\dim(V) = n$, $\dim(W) = m$. Choose a basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ of V and a basis $\mathbf{f}_1, \dots, \mathbf{f}_m$ of W . Note that formally what we are doing here is choosing *ordered* bases- above we defined a basis of a vector space to be simply a *subset* without any preferred ordering, but here we actually make a choice of two ordered sets of bases, $\mathbf{E} = (\mathbf{e}_1, \dots, \mathbf{e}_n)$ and $\mathbf{F} = (\mathbf{f}_1, \dots, \mathbf{f}_m)$, the ordering being encoded in the choice of indices.

Now, for $1 \leq j \leq n$, $T(\mathbf{e}_j) \in W$, so $T(\mathbf{e}_j)$ can be written uniquely as a linear combination of $\mathbf{f}_1, \dots, \mathbf{f}_m$. Let

$$\begin{aligned} T(\mathbf{e}_1) &= \alpha_{11}\mathbf{f}_1 + \alpha_{21}\mathbf{f}_2 + \cdots + \alpha_{m1}\mathbf{f}_m \\ T(\mathbf{e}_2) &= \alpha_{12}\mathbf{f}_1 + \alpha_{22}\mathbf{f}_2 + \cdots + \alpha_{m2}\mathbf{f}_m \\ &\vdots \\ T(\mathbf{e}_n) &= \alpha_{1n}\mathbf{f}_1 + \alpha_{2n}\mathbf{f}_2 + \cdots + \alpha_{mn}\mathbf{f}_m \end{aligned}$$

where the coefficients $\alpha_{ij} \in K$ (for $1 \leq i \leq m$, $1 \leq j \leq n$) are uniquely determined.

The coefficients α_{ij} form an $m \times n$ matrix

$$A = \begin{pmatrix} \alpha_{11} & \alpha_{12} & \cdots & \alpha_{1n} \\ \alpha_{21} & \alpha_{22} & \cdots & \alpha_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_{m1} & \alpha_{m2} & \cdots & \alpha_{mn} \end{pmatrix}$$

over K . Then A is called the matrix of the linear map T with respect to the chosen bases of V and W . Note that the columns of A are the images $T(\mathbf{e}_1), \dots, T(\mathbf{e}_n)$ of the basis vectors of V represented as column vectors with respect to the basis $\mathbf{f}_1, \dots, \mathbf{f}_m$ of W .

It was shown in MA106 that T is uniquely determined by A , and so there is a one-one correspondence between linear maps $T : V \rightarrow W$ and $m \times n$ matrices over K , which depends on the choice of ordered bases of V and W .

1 Review of some MA106 material

For $\mathbf{v} \in V$, we can write \mathbf{v} uniquely as a linear combination of the basis vectors \mathbf{e}_i ; that is, $\mathbf{v} = x_1\mathbf{e}_1 + \cdots + x_n\mathbf{e}_n$, where the x_i are uniquely determined by \mathbf{v} and the basis \mathbf{e}_i . We shall call x_i the *coordinates* of \mathbf{v} with respect to the basis $\mathbf{e}_1, \dots, \mathbf{e}_n$. We associate the column vector

$$\underline{\mathbf{v}} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in K^{n,1},$$

to \mathbf{v} , where $K^{n,1}$ denotes the space of $n \times 1$ -column vectors with entries in K . Notice that $\underline{\mathbf{v}}$ depends on the chosen basis \mathbf{E} so a notation such as $\mathbf{v}_{\mathbf{E}}$ or $\underline{\mathbf{v}}_{\mathbf{E}}$ would possibly be better, but also heavier, so we stick with $\underline{\mathbf{v}}$ and assume you bear in mind that $\underline{\mathbf{v}}$ not only depends on \mathbf{v} but also on \mathbf{E} .

It was proved in MA106 that if A is the matrix of the linear map T , then for $\mathbf{v} \in V$, we have $T(\mathbf{v}) = \mathbf{w}$ if and only if $A\underline{\mathbf{v}} = \underline{\mathbf{w}}$, where $\underline{\mathbf{w}} \in K^{m,1}$ is the column vector associated with $\mathbf{w} \in W$.

Example. We can write down the matrices for the linear maps in Example 2, using the standard bases for V and W : the standard basis of \mathbb{R}^n is e_1, \dots, e_n where e_i is the column vector with a 1 in the i th row and all other entries 0 (so it's the $n \times 1$ matrix defined by $\alpha_{j,1} = 1$ if $j = i$ and $\alpha_{j,i} = 0$ otherwise).

1. We calculate that $T(e_1) = e_1 = 1 \cdot e_1 + 0 \cdot e_2$, $T(e_2) = e_2 = 0 \cdot e_1 + 1 \cdot e_2$ and $T(e_3) = 0 = 0 \cdot e_1 + 0 \cdot e_2$ (OK, this could be confusing so we could denote the standard basis for W by f_1, f_2). The matrix is thus

$$A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}.$$

2. We skip the details but the matrix is

$$A = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}.$$

3. This time $T(e_1) = e_1$, $T(e_2) = e_1 + e_2$ and $T(e_3) = e_2$ and so the matrix is

$$A = \begin{pmatrix} 1 & 1 & 0 \\ 0 & 1 & 1 \end{pmatrix}.$$

1.5 Change of basis

Let V be a vector space of dimension n over a field K , and let $\mathbf{e}_1, \dots, \mathbf{e}_n$ and $\mathbf{e}'_1, \dots, \mathbf{e}'_n$ be two bases of V (ordered of course). Then there is an invertible $n \times n$ matrix $P = (p_{ij})$ such that

$$\mathbf{e}'_j = \sum_{i=1}^n p_{ij}\mathbf{e}_i \quad \text{for } 1 \leq j \leq n. \quad (*)$$

1 Review of some MA106 material

Note that the columns of P are the new basis vectors \mathbf{e}'_i written as column vectors in the old basis vectors \mathbf{e}_i . (Recall also that P is the matrix of the identity map $V \rightarrow V$ using basis $\mathbf{e}'_1, \dots, \mathbf{e}'_n$ in the domain and basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ in the codomain.)

Often, but not always, the original basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ will be the standard basis of K^n .

Example. Let $V = \mathbb{R}^3$, $\mathbf{e}_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$, $\mathbf{e}_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$, $\mathbf{e}_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$ (the standard basis) and $\mathbf{e}'_1 = \begin{pmatrix} 0 \\ 1 \\ 2 \end{pmatrix}$, $\mathbf{e}'_2 = \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}$, $\mathbf{e}'_3 = \begin{pmatrix} -1 \\ 0 \\ 0 \end{pmatrix}$. Then

$$P = \begin{pmatrix} 0 & 1 & -1 \\ 1 & 2 & 0 \\ 2 & 0 & 0 \end{pmatrix}.$$

The following result was proved in MA106.

Proposition 1.5.1. *With the above notation, let $\mathbf{v} \in V$, and let $\underline{\mathbf{v}}$ and $\underline{\mathbf{v}}'$ denote the column vectors associated with \mathbf{v} when we use the bases $\mathbf{e}_1, \dots, \mathbf{e}_n$ and $\mathbf{e}'_1, \dots, \mathbf{e}'_n$, respectively. Then $P\underline{\mathbf{v}}' = \underline{\mathbf{v}}$.*

So, in the example above, if we take $\mathbf{v} = \begin{pmatrix} 1 \\ -2 \\ 4 \end{pmatrix}$, then we have $\mathbf{v} = \mathbf{e}_1 - 2\mathbf{e}_2 + 4\mathbf{e}_3$ (obviously);

so the coordinates of \mathbf{v} in the basis $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$ are $\underline{\mathbf{v}} = \begin{pmatrix} 1 \\ -2 \\ 4 \end{pmatrix}$.

On the other hand, we also have $\mathbf{v} = 2\mathbf{e}'_1 - 2\mathbf{e}'_2 - 3\mathbf{e}'_3$, so the coordinates of \mathbf{v} in the basis $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3\}$ are

$$\underline{\mathbf{v}}' = \begin{pmatrix} 2 \\ -2 \\ -3 \end{pmatrix},$$

and you can check that

$$P\underline{\mathbf{v}}' = \begin{pmatrix} 0 & 1 & -1 \\ 1 & 2 & 0 \\ 2 & 0 & 0 \end{pmatrix} \begin{pmatrix} 2 \\ -2 \\ -3 \end{pmatrix} = \begin{pmatrix} 1 \\ -2 \\ 4 \end{pmatrix} = \underline{\mathbf{v}},$$

just as Proposition 1.5.1 says.

This equation $P\underline{\mathbf{v}}' = \underline{\mathbf{v}}$ describes the change of coordinates associated with our basis change. If we want to compute the new coordinates from the old ones, we need to use the inverse matrix: $\underline{\mathbf{v}}' = P^{-1}\underline{\mathbf{v}}$. Thus, to enable calculations in the new basis we need both matrices P and P^{-1} . We'll be using this relationship over and over again, so make sure you're happy with it!

Which matrix, P or P^{-1} should be called the *basis change matrix* or *transition matrix* from the original basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ to the new basis $\mathbf{e}'_1, \dots, \mathbf{e}'_n$?

1 Review of some MA106 material

Well, the books are split on this. As a historic quirk, the *basis change matrix* in Algebra-1 was always P and the *basis change matrix* in Linear Algebra was P^{-1} since around 2011. We continue with this noble tradition of calling P the *basis change matrix* because, otherwise, we risk introducing typos throughout the text.

Now let $T : V \rightarrow W$, $\mathbf{e}_i, \mathbf{f}_i$ and A be as in Subsection 1.4 above, and choose new bases $\mathbf{e}'_1, \dots, \mathbf{e}'_n$ of V and $\mathbf{f}'_1, \dots, \mathbf{f}'_m$ of W . Then

$$T(\mathbf{e}'_j) = \sum_{i=1}^m \beta_{ij} \mathbf{f}'_i \text{ for } 1 \leq j \leq n,$$

where $B = (\beta_{ij})$ is the $m \times n$ matrix of T with respect to the bases $\{\mathbf{e}'_i\}$ and $\{\mathbf{f}'_i\}$ of V and W . Let the $n \times n$ matrix $P = (p_{ij})$ be the basis change matrix for the original basis $\{\mathbf{e}_i\}$ and new basis $\{\mathbf{e}'_i\}$, and let the $m \times m$ matrix $Q = (q_{ij})$ be the basis change matrix for original basis $\{\mathbf{f}_i\}$ and new basis $\{\mathbf{f}'_i\}$. The following theorem was proved in MA106:

Theorem 1.5.2. *With the above notation, we have $AP = QB$, or equivalently $B = Q^{-1}AP$.*

In most of the applications in this module we will have $V = W (= K^n)$, $\{\mathbf{e}_i\} = \{\mathbf{f}_i\}$, and $\{\mathbf{e}'_i\} = \{\mathbf{f}'_i\}$. So $P = Q$, and hence $B = P^{-1}AP$.

You may have noticed that the above is a bit messy, and it can be difficult to remember the definitions of P and Q (and to distinguish them from their inverses). Experience shows that students (and lecturers) have trouble with this. So here is what I hope is a better and more transparent way to think about change of basis in vector spaces and the way it affects representing matrices for linear maps:

First, we saw in the preceding section, that given:

1. a linear map $T : V \rightarrow W$, $\dim(V) = n$, $\dim(W) = m$;
2. ordered bases $\mathbf{E} = (\mathbf{e}_1, \dots, \mathbf{e}_n)$ and $\mathbf{F} = (\mathbf{f}_1, \dots, \mathbf{f}_m)$ of V and W ;

we can associate to T an $m \times n$ -matrix in $K^{m \times n}$ representing the linear map T with respect to the chosen ordered bases. Let's do our book-keeping neatly and try to keep track of all the data involved in our notation: let's denote this matrix temporarily by

$$\mathcal{M}(T)_{\mathbf{E}}^{\mathbf{F}}$$

Note that the lower index \mathbf{E} remembers the basis in the source V , the upper index \mathbf{F} remembers the basis in the target, and \mathcal{M} just stands for matrix. Of course that's a notational monstrosity, but you will see that for the purpose of explaining base change, it is very convenient. Indeed, choosing different ordered bases for V and W ,

$$\mathbf{E}' = (\mathbf{e}'_1, \dots, \mathbf{e}'_n) \text{ and } \mathbf{F}' = (\mathbf{f}'_1, \dots, \mathbf{f}'_m)$$

the problem we want to address is: how are the matrices

$$A = \mathcal{M}(T)_{\mathbf{E}}^{\mathbf{F}} \text{ and } B = \mathcal{M}(T)_{\mathbf{E}'}^{\mathbf{F}'}$$

1 Review of some MA106 material

related? The answer to this is very easy if you remember from MA106 that matrix multiplication is compatible with composition of linear maps in the following sense: suppose

$$\begin{array}{ccccc} U & \xrightarrow{R} & V & \xrightarrow{S} & W \\ \mathbf{A} & & \mathbf{B} & & \mathbf{C} \end{array}$$

is a diagram of vector spaces and linear maps, and $\mathbf{A}, \mathbf{B}, \mathbf{C}$ are ordered bases in U, V, W . Then we have *the very basic fact* that

$$\mathcal{M}_{\mathbf{A}}^{\mathbf{C}}(S \circ R) = \mathcal{M}_{\mathbf{B}}^{\mathbf{C}}(S) \cdot \mathcal{M}_{\mathbf{A}}^{\mathbf{B}}(R).$$

Don't be intimidated by the formula and take a second to think about how natural this is! If we form the composite map $S \circ R$ and pass to the matrix representing it with respect to the given ordered bases, we can also get it by matrix-multiplying the matrices for S and R with respect to the chosen ordered bases! Now back to our problem above: consider the sequence of linear maps between vector spaces together with choices of ordered bases:

$$\begin{array}{ccccccc} V & \xrightarrow{\text{id}_V} & V & \xrightarrow{T} & W & \xrightarrow{\text{id}_W} & W \\ \mathbf{E}' & & \mathbf{E} & & \mathbf{F} & & \mathbf{F}' \end{array}$$

Applying the preceding basic fact gives

$$\mathcal{M}(T)_{\mathbf{E}'}^{\mathbf{F}'} = \mathcal{M}(\text{id}_W)_{\mathbf{F}}^{\mathbf{F}'} \cdot \mathcal{M}(T)_{\mathbf{E}}^{\mathbf{F}} \cdot \mathcal{M}(\text{id}_V)_{\mathbf{E}'}^{\mathbf{E}}.$$

Or, putting

$$P := \mathcal{M}(\text{id}_V)_{\mathbf{E}'}^{\mathbf{E}}, \quad Q := \mathcal{M}(\text{id}_W)_{\mathbf{F}}^{\mathbf{F}'}$$

and noticing that

$$\mathcal{M}(\text{id}_W)_{\mathbf{F}}^{\mathbf{F}'} = (\mathcal{M}(\text{id}_W)_{\mathbf{F}'}^{\mathbf{F}})^{-1}$$

we get

$$B = Q^{-1}AP$$

which *proves* Theorem 1.5.2, *but* also gives us a means to *remember* the right definitions of P and Q (which is important because that is the vital information and this is precisely the information students and lecturer always tend to forget): for example, $P = \mathcal{M}(\text{id}_V)_{\mathbf{E}'}^{\mathbf{E}}$ is the matrix whose columns are the basis vectors \mathbf{e}'_i written in the old basis \mathbf{E} with basis vectors \mathbf{e}_i . You don't have to remember the entire discussion preceding Theorem 1.5.2 anymore (which is necessary to understand what the theorem says): it's all encoded in the notation! I hope you will never forget this base change formula again.

2 The Jordan Canonical Form

2.1 Introduction

Throughout this section V will be a vector space of dimension n over a field K , $T : V \rightarrow V$ will be a linear map, and A will be the matrix of T with respect to a fixed basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ of V (the same in the source and target V). Our aim is to find a new basis $\mathbf{e}'_1, \dots, \mathbf{e}'_n$ for V , such that the matrix of T with respect to the new basis is as simple as possible. Equivalently (by Theorem 1.5.2), we want to find an invertible matrix P (the associated basis change matrix) such that $P^{-1}AP$ is as simple as possible.

(Recall that if B is a matrix which can be written in the form $B = P^{-1}AP$, we say B is *similar* to A . So a third way of saying the above is that we want to find a matrix that's similar to A , but which is as nice as possible.)

One particularly simple form of a matrix is a diagonal matrix. So we'd really rather like it if every matrix was similar to a diagonal matrix. But this won't work: we saw in MA106 that the matrix $\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$, for example, is not similar to a diagonal matrix. (We say this matrix is not *diagonalizable*.)

The point of this section of the module is to show that although we can't always get to a diagonal matrix, we can get pretty close (at least if K is \mathbb{C}). Under this assumption, it can be proved that A is always similar to a matrix B of a certain type (called the *Jordan canonical form* or sometimes *Jordan normal form* of the matrix), which is not far off being diagonal: its only non-zero entries are on the diagonal or just above it.

2.2 Eigenvalues and eigenvectors

We start by summarising some of what we know from MA106 which is going to be relevant to us here.

If we can find some $0 \neq \mathbf{v} \in V$ and $\lambda \in K$ such that $T\mathbf{v} = \lambda\mathbf{v}$, or equivalently $A\mathbf{v} = \lambda\mathbf{v}$, then λ is an *eigenvalue*, and \mathbf{v} a corresponding *eigenvector* of T (or of A).

From MA106, you have a theorem that tells you when a matrix is diagonalizable:

Proposition 2.2.1. *Let $T : V \rightarrow V$ be a linear map. Then the matrix of T is diagonal with respect to some basis of V if and only if V has a basis consisting of eigenvectors of T .*

This is a nice theorem, but it is also more or less a tautology, and it doesn't tell you how you might find such a basis! But there's one case where it's easy, as another theorem from MA106 tells us:

Proposition 2.2.2. *Let $\lambda_1, \dots, \lambda_r$ be distinct eigenvalues of $T : V \rightarrow V$, and let $\mathbf{v}_1, \dots, \mathbf{v}_r$ be corresponding eigenvectors. (So $T(\mathbf{v}_i) = \lambda_i\mathbf{v}_i$ for $1 \leq i \leq r$.) Then $\mathbf{v}_1, \dots, \mathbf{v}_r$ are linearly independent.*

2 The Jordan Canonical Form

Corollary 2.2.3. *If the linear map $T : V \rightarrow V$ (or equivalently the $n \times n$ matrix A) has n distinct eigenvalues, where $n = \dim(V)$, then T (or A) is diagonalizable.*

2.3 The minimal polynomial

The minimal polynomial, while arguably not the most important player in the spectral theory of endomorphisms, derives its importance from the fact that it can be used to detect diagonalisability and also classifies nilpotent transformations, and we'll start with it to get off the ground.

If $A \in K^{n,n}$ is a square $n \times n$ matrix over K , and $p \in K[x]$ is a polynomial, then we can make sense of $p(A)$: we just calculate the powers of A in the usual way, and then plug them into the formula defining p , interpreting the constant term as a multiple of I_n .

For instance, if $K = \mathbb{Q}$, $p = 2x^2 - \frac{3}{2}x + 11$, and $A = \begin{pmatrix} 2 & 3 \\ 0 & 1 \end{pmatrix}$, then $A^2 = \begin{pmatrix} 4 & 9 \\ 0 & 1 \end{pmatrix}$, and

$$\begin{aligned} p(A) &= 2 \begin{pmatrix} 4 & 9 \\ 0 & 1 \end{pmatrix} - \frac{3}{2} \begin{pmatrix} 2 & 3 \\ 0 & 1 \end{pmatrix} + 11 \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \\ &= \begin{pmatrix} 16 & 27/2 \\ 0 & 23/2 \end{pmatrix}. \end{aligned}$$

Warning. Notice that this is in general of course *not* the same as the matrix $\begin{pmatrix} p(2) & p(3) \\ p(0) & p(1) \end{pmatrix}$.

Theorem 2.3.1. *Let $A \in K^{n,n}$. Then there is some non-zero polynomial $p \in K[x]$ of degree at most n^2 such that $p(A)$ is the $n \times n$ zero matrix $\mathbf{0}_n$.*

Proof. The key thing to observe is that $K^{n,n}$, the space of $n \times n$ matrices over K , is itself a vector space over K . Its dimension is n^2 .

Let's consider the set $\{I_n, A, A^2, \dots, A^{n^2}\} \subset K^{n,n}$. Since this is a set of $n^2 + 1$ vectors in an n^2 -dimensional vector space, there is a nontrivial linear dependency relation between them. That is, we can find constants $\lambda_0, \lambda_1, \dots, \lambda_{n^2}$, not all zero, such that

$$\lambda_0 I_n + \dots + \lambda_{n^2} A^{n^2} = \mathbf{0}_n.$$

Now we define the polynomial $p = \lambda_0 + \lambda_1 x + \dots + \lambda_{n^2} x^{n^2}$. This isn't zero, and its degree is at most n^2 . (It might be less, since λ_{n^2} might be 0.) Then that's it! □

Is there a way of finding a unique polynomial (of minimal degree) that A satisfies? To answer that question, we'll have to think a little bit about arithmetic in $K[x]$.

Note that we can do "division" with polynomials, a bit like with integers. We can divide one polynomial p (with $p \neq 0$) into another polynomial q and get a remainder with degree less than

2 The Jordan Canonical Form

p . For example, if $q = x^5 - 3$, $p = x^2 + x + 1$, then we find $q = sp + r$ with $s = x^3 - x^2 + 1$ and $r = -x - 4$.

If the remainder is 0, so $q = sp$ for some s , we say “ p divides q ” and write this relation as $p \mid q$.

Finally, a polynomial with coefficients in a field K is called *monic* if the coefficient of the highest power of x is 1. So, for example, $x^3 - 2x^2 + x + 11$ is monic, but $2x^2 - x - 1$ is not.

Theorem 2.3.2. *Let A be an $n \times n$ matrix over K representing the linear map $T : V \rightarrow V$. Then*

- (i) *There is a unique monic non-zero polynomial $p(x)$ with minimal degree and coefficients in K such that $p(A) = 0$.*
- (ii) *If $q(x)$ is any polynomial with $q(A) = 0$, then $p \mid q$.*

Proof. (i) If we have any polynomial $p(x)$ with $p(A) = 0$, then we can make p monic by multiplying it by a constant. By Theorem 2.3.1, there exists such a $p(x)$, so there exists one of minimal degree. If we had two distinct monic polynomials $p_1(x)$, $p_2(x)$ of the same minimal degree with $p_1(A) = p_2(A) = 0$, then $p = p_1 - p_2$ would be a non-zero polynomial of smaller degree with $p(A) = 0$, contradicting the minimality of the degree, so p is unique.

(ii) Let $p(x)$ be the minimal monic polynomial in (i) and suppose that $q(A) = 0$. As we saw above, we can write $q = sp + r$ where r has smaller degree than p . If r is non-zero, then $r(A) = q(A) - s(A)p(A) = 0$ contradicting the minimality of p , so $r = 0$ and $p \mid q$. \square

Definition 2.3.3. The unique monic non-zero polynomial $\mu_A(x)$ of minimal degree with $\mu_A(A) = 0$ is called the *minimal polynomial* of A .

We know that for $p \in K[x]$, $p(T) = \mathbf{0}_V$ if and only if $p(A) = \mathbf{0}_n$; so μ_A is also the unique monic polynomial of minimal degree such that $\mu_A(T) = 0$ (the minimal polynomial of T .) In particular, since similar matrices A and B represent the same linear map T , and their minimal polynomial is the same as that of T , we have

Proposition 2.3.4. *Similar matrices have the same minimal polynomial.*

By Theorem 2.3.1 and Theorem 2.3.2 (ii), we have

Corollary 2.3.5. *The minimal polynomial of an $n \times n$ matrix A has degree at most n^2 .*

(In the next section, we’ll see that we can do much better than this.)

Example. If D is a diagonal matrix, say

$$D = \begin{pmatrix} d_{11} & & \\ & \ddots & \\ & & d_{nn} \end{pmatrix},$$

then for any polynomial p we see that $p(D)$ is the diagonal matrix with entries

$$\begin{pmatrix} p(d_{11}) & & \\ & \ddots & \\ & & p(d_{nn}) \end{pmatrix}.$$

2 The Jordan Canonical Form

Hence $p(D) = 0$ if and only if $p(d_{ii}) = 0$ for all i . So for instance if

$$D = \begin{pmatrix} 3 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 2 \end{pmatrix},$$

the minimal polynomial of D is the smallest-degree polynomial which has 2 and 3 as roots, which is clearly $\mu_D(x) = (x - 2)(x - 3) = x^2 - 5x + 6$.

We can generalize this example as follows

Proposition 2.3.6. *Let D be any diagonal matrix and let $\{\delta_1, \dots, \delta_r\}$ be the set of diagonal entries of D (i.e. without any repetitions, so the values $\delta_1, \dots, \delta_r$ are all different). Then we have*

$$\mu_D(x) = (x - \delta_1)(x - \delta_2) \dots (x - \delta_r).$$

Proof. As in the example, we have $p(D) = 0$ if and only if $p(\delta_i) = 0$ for all $i \in \{1, \dots, r\}$. The smallest-degree monic polynomial vanishing at these points is clearly the polynomial above. \square

Corollary 2.3.7. *If A is any diagonalizable matrix, then $\mu_A(x)$ is a product of distinct linear factors.*

Proof. Clear from Proposition 2.3.6 and Proposition 2.3.4. \square

Remark. *We'll see later in the course that this is a necessary and sufficient condition: A is diagonalizable if and only if $\mu_A(x)$ is a product of distinct linear factors. But we don't have enough tools to prove this theorem yet – be patient!*

2.4 The Cayley–Hamilton theorem

Theorem 2.4.1 (Cayley–Hamilton). *Let $c_A(x)$ be the characteristic polynomial of the $n \times n$ matrix A over an arbitrary field K . Then $c_A(A) = \mathbf{0}$.*

Proof. Let's agree to drop the various subscripts and bold zeroes – it'll be obvious from context when we mean a zero matrix, zero vector, zero linear map, etc.

Recall from MA106 that, if B is any $n \times n$ matrix, the “adjugate matrix” of B is another matrix $\text{adj}(B)$ which was constructed along the way to constructing the inverse of B . The entries of $\text{adj}(B)$ are the “cofactors” of B : the (i, j) entry of B is $(-1)^{i+j}c_{ji}$ (note the transposition of indices here!), where $c_{ji} = \det(B_{ji})$, B_{ji} being the $(n - 1) \times (n - 1)$ matrix obtained by deleting the j -th row and the i -th column of B . The key property of $\text{adj}(B)$ is that it satisfies

$$B \text{adj}(B) = \text{adj}(B)B = (\det B)I_n.$$

(Notice that if B is invertible, this just says that $\text{adj}(B) = (\det B)B^{-1}$, but the adjugate matrix still makes sense even if B is not invertible.)

2 The Jordan Canonical Form

Let's apply this to the matrix $B = A - xI_n$. By definition, $\det(B)$ is the characteristic polynomial $c_A(x)$, so

$$\text{adj}(A - xI_n)(A - xI_n) = c_A(x)I_n. \quad (2)$$

Now we use the following statement whose proof is obvious: suppose $P(x) = \sum_j P_j x^j$ and $Q(x) = \sum_k Q_k x^k$ are two polynomials in the indeterminate x with *matrix coefficients*; so P_j and Q_k are $n \times n$ matrices. Then the product of P and Q is $R(x) = \sum_l R_l x^l$ with

$$R_l = \sum_{j+k=l} P_j Q_k.$$

Then if an $n \times n$ matrix M commutes with all the coefficients of Q we have $R(M) = P(M)Q(M)$. We now apply this observation with

$$P(x) = \text{adj}(A - xI_n), \quad Q(x) = A - xI_n, \quad M = A.$$

Since $Q(A) = 0$, we get $c_A(A) = 0$. □

Corollary 2.4.2. *For any $A \in K^{n,n}$, we have $\mu_A \mid c_A$, and in particular $\deg(\mu_A) \leq n$.*

Example. Let D be the diagonal matrix $\begin{pmatrix} 3 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 2 \end{pmatrix}$ from the previous example. We saw above that $\mu_A(x) = (x - 2)(x - 3)$. However, it's easy to see that

$$c_A(x) = \begin{vmatrix} 3-x & 0 & 0 \\ 0 & 3-x & 0 \\ 0 & 0 & 2-x \end{vmatrix} = -(x-2)(x-3)^2.$$

How NOT to prove the Cayley–Hamilton theorem It is very tempting to try and prove the Cayley–Hamilton theorem as follows: we know that

$$c_A(x) = \det(A - xI_n),$$

so shouldn't we have

$$c_A(A) = \det(A - AI_n) = \det(A - A) = \det(0) = 0?$$

This is **wrong**. In fact, $c_A(A)$ is a matrix, and $\det(A - AI_n)$ is an element of K , so they are not even objects of the same type in general.

2.5 Calculating the minimal polynomial

We will present two methods for this.

Method 1 (“top down”; always never works in practice; it only works well if a benign lecturer or some other benevolent power reveals to you the factorisation of the characteristic polynomial into irreducibles).

2 The Jordan Canonical Form

Lemma 2.5.1. *Let λ be any eigenvalue of A . Then $\mu_A(\lambda) = 0$.*

Proof. Let $\underline{v} \in K^{n,1}$ be an eigenvector corresponding to λ . Then $A^n \underline{v} = \lambda^n \underline{v}$, and hence for any polynomial $p \in K[x]$, we have

$$p(A)\underline{v} = p(\lambda)\underline{v}.$$

We know that $\mu_A(A)\underline{v} = 0$, since $\mu_A(A)$ is the zero matrix. Hence $\mu_A(\lambda)\underline{v} = 0$, and since $\underline{v} \neq 0$ and $\mu_A(\lambda)$ is an element of K (not a matrix!), this can only happen if $\mu_A(\lambda) = 0$. \square

This lemma, together with Cayley–Hamilton, give us very, very few possibilities for μ_A . Let's look at an example.

Example. Take $K = \mathbb{C}$ and let

$$A = \begin{pmatrix} 4 & 0 & -1 & -1 \\ 1 & 2 & 0 & 0 \\ 2 & -2 & 2 & -2 \\ -1 & 1 & 0 & 3 \end{pmatrix}.$$

This is rather large, but it has a fair few zeros, so you can calculate its characteristic polynomial fairly quickly by hand and find out that

$$c_A(x) = x^4 - 11x^3 + 45x^2 - 81x + 54.$$

Some trial and error shows that 2 is a root of this, and we find that

$$c_A(x) = (x - 2)(x^3 - 9x^2 + 27x - 27) = (x - 2)(x - 3)^3.$$

So $\mu_A(x)$ divides $(x - 2)(x - 3)^3$. On the other hand, the eigenvalues of A are the roots of $c_A(x)$, namely $\{2, 3\}$; and we know that μ_A must have each of these as roots. So the only possibilities for $\mu_A(x)$ are:

$$\mu_A(x) \in \left\{ \begin{array}{l} (x - 2)(x - 3), \\ (x - 2)(x - 3)^2, \\ (x - 2)(x - 3)^3. \end{array} \right\}.$$

Some slightly tedious calculation shows that $(A - 2)(A - 3)$ isn't zero, and nor is $(A - 2)(A - 3)^2$, and so it must be the case that $(x - 2)(x - 3)^3$ is the minimal polynomial of A .

Method 2 (“bottom up”; this works well, also for large matrices)

This is based on

Lemma 2.5.2. *Let $T: V \rightarrow V$ be a linear map of an n -dimensional vector space V over K to itself, and suppose W_1, \dots, W_k are finitely many T -invariant subspaces spanning V . In other words, we require $T(W_i) \subset W_i$ and*

$$V = W_1 + \dots + W_k$$

2 The Jordan Canonical Form

(but the sum doesn't have to be direct). Let $\mu_i(x)$ be the minimal polynomial of $T|_{W_i}$. Then

$$\mu_T(x) = \text{l.c.m.}\{\mu_1, \dots, \mu_k\}.$$

In words: the minimal polynomial of T is the least common multiple of the minimal polynomials of the $T|_{W_i}$, $i = 1, \dots, k$.

Proof. First we will show that setting

$$f(x) = \text{l.c.m.}\{\mu_1, \dots, \mu_k\}$$

we have that $\mu_T(x)$ divides $f(x)$. Indeed, if $v \in W_i$, then writing $f(x) = g_i(x)\mu_i(x)$ we calculate

$$f(T)v = g_i(T)\mu_i(T)v = g_i(T|_{W_i})\mu_i(T|_{W_i})v = 0$$

since $\mu_i(T|_{W_i}) = 0$. Since this argument is valid for any i and the W_i 's span V , we conclude that $f(T)$ annihilates all of V hence is the zero linear map on V . Thus $f(x)$ is divisible by $\mu_T(x)$.

But $f(x)$ also divides $\mu_T(x)$: indeed, $\mu_T(T) = 0$, and hence also $\mu_T(T|_{W_i}) = 0$ for any i . Hence, $\mu_T(x)$ is divisible by any $\mu_i(x)$, and consequently by their least common multiple, too.

Since both $f(x)$ and $\mu_T(x)$ are monic, they must be equal. □

The preceding Lemma allows us to come up with a sensible algorithm to compute the minimal polynomial of T :

Algorithm:

Pick any $v \neq 0$ in V and set

$$W = \text{span}\{v, T(v), T^2(v), \dots\}.$$

By definition, W is T -invariant. Now let d be the minimal positive integer such that

$$v, T(v), \dots, T^d(v)$$

are linearly dependent. In particular,

$$v, T(v), \dots, T^{d-1}(v)$$

are linearly independent, and if $p(x)$ is any polynomial of degree $\leq d-1$, $p(T)v$ will never be zero: hence the minimal polynomial $\mu_{T|_W}(x)$ has degree $\geq d$. There is a nontrivial linear dependency relation of the form

$$T^d(v) + c_{d-1}T^{d-1}(v) + \dots + c_1T(v) + c_0v = 0.$$

Consider the polynomial

$$x^d + c_{d-1}x^{d-1} + \dots + c_1x + c_0.$$

2 The Jordan Canonical Form

We claim this must be $\mu_{T|_W}(x)$: indeed, it is monic, $\mu_{T|_W}(T|_W)$ annihilates W , and $\mu_{T|_W}(x)$ is of smallest possible degree d with this property. Therefore we have computed $\mu_{T|_W}(x)$, and we can set

$$W_1 := W, \mu_1(x) := \mu_{T|_W}(x).$$

If $W_1 \neq V$, pick a vector v' not in W_1 and repeat the preceding procedure, leading to a T -invariant subspace W_2 such that the span of W_1 and W_2 will be strictly larger than W_1 . Since V is finite-dimensional, after finitely many steps, we compute in this way W_1, \dots, W_k and polynomials $\mu_1(x), \dots, \mu_k(x)$ satisfying the conditions in Lemma 2.5.2. Since computing a least common multiple presents no problem (use the Euclidean algorithm for polynomials repeatedly), we are done.

2.6 Jordan chains and Jordan blocks

We'll now consider some special vectors attached to our matrix, which satisfy a condition a little like eigenvectors (but weaker). These will be the stepping-stones towards the Jordan canonical form.

Definition 2.6.1. A non-zero vector $\mathbf{v} \in K^{n,1}$ such that $(A - \lambda I_n)^i \mathbf{v} = \mathbf{0}$, for some $i > 0$, is called a *generalised eigenvector* of A with respect to the eigenvalue λ .

Note that, for fixed $i > 0$,

$$N_i(A, \lambda) := \{ \mathbf{v} \in V \mid (A - \lambda I_n)^i \mathbf{v} = \mathbf{0} \}$$

is the nullspace of $(A - \lambda I_n)^i$, and is called the *generalised eigenspace of index i* of A with respect to λ .

The generalised eigenspace of index 1 is just called the *eigenspace* of A w.r.t. λ ; it consists of the eigenvectors of A w.r.t. λ , together with the zero vector. We sometimes also consider the *full generalised eigenspace* of A w.r.t. λ , which is the set of all generalised eigenvectors together with the zero vector; this is the union of the generalised eigenspaces of index i over all $i \in \mathbb{N}$.

We can arrange generalised eigenvectors into “chains”:

Definition 2.6.2. A *Jordan chain of length k* is a sequence of non-zero vectors $\mathbf{v}_1, \dots, \mathbf{v}_k \in K^{n,1}$ that satisfies

$$A\mathbf{v}_1 = \lambda\mathbf{v}_1, \quad A\mathbf{v}_i = \lambda\mathbf{v}_i + \mathbf{v}_{i-1}, \quad 2 \leq i \leq k,$$

for some eigenvalue λ of A .

Equivalently, $(A - \lambda I_n)\mathbf{v}_1 = \mathbf{0}$ and $(A - \lambda I_n)\mathbf{v}_i = \mathbf{v}_{i-1}$ for $2 \leq i \leq k$, so $(A - \lambda I_n)^i \mathbf{v}_i = \mathbf{0}$ for $1 \leq i \leq k$. Thus all of the vectors in a Jordan chain are generalised eigenvectors, and \mathbf{v}_i lies in the generalised eigenspace of index i .

Lemma 2.6.3. Let $\mathbf{v}_1, \dots, \mathbf{v}_k \in K^{n,1}$ be a Jordan chain of length k for eigenvalue λ of A . Then $\mathbf{v}_1, \dots, \mathbf{v}_k$ are linearly independent.

2 The Jordan Canonical Form

Proof. Exercise. □

For example, take $K = \mathbb{C}$ and consider the matrix

$$A = \begin{pmatrix} 3 & 1 & 0 \\ 0 & 3 & 1 \\ 0 & 0 & 3 \end{pmatrix}.$$

We see that, for $\{\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3\}$ the standard basis of $\mathbb{C}^{3,1}$, we have $A\mathbf{b}_1 = 3\mathbf{b}_1$, $A\mathbf{b}_2 = 3\mathbf{b}_2 + \mathbf{b}_1$, $A\mathbf{b}_3 = 3\mathbf{b}_3 + \mathbf{b}_2$, so $\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3$ is a Jordan chain of length 3 for the eigenvalue 3 of A . The generalised eigenspaces of index 1, 2, and 3 are respectively $\langle \mathbf{b}_1 \rangle$, $\langle \mathbf{b}_1, \mathbf{b}_2 \rangle$, and $\langle \mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3 \rangle$.

Note that this isn't the only possible Jordan chain. Obviously, $\{17\mathbf{b}_1, 17\mathbf{b}_2, 17\mathbf{b}_3\}$ would be a Jordan chain; but there are more devious possibilities – you can check that $\{\mathbf{b}_1, \mathbf{b}_1 + \mathbf{b}_2, \mathbf{b}_2 + \mathbf{b}_3\}$ is a Jordan chain, so there can be several Jordan chains with the same first vector. On the other hand, two Jordan chains with the same *last* vector are the same and in particular have the same length.

What are the generalised eigenspaces here? The only eigenvalue is 3. For this eigenvalue, the generalised eigenspace of index 1 is $\langle \mathbf{b}_1 \rangle$ (the linear span of \mathbf{b}_1); the generalised eigenspace of index 2 is $\langle \mathbf{b}_1, \mathbf{b}_2 \rangle$; and the generalised eigenspace of index 3 is the whole space $\langle \mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3 \rangle$. So the dimensions are $(1, 2, 3)$.

Proposition 2.6.4. *The dimensions of corresponding generalised eigenspaces of similar matrices are the same.*

Proof. Notice that the dimension of a generalised eigenspace of A is the nullity of $(T - \lambda I_V)^i$, which depends only on the linear map T associated with A . Therefore it's independent of the choice of basis. Since similar matrices represent the same linear map in different bases, the proposition follows. □

In the example we had earlier, the standard basis of $K^{n,1}$ was a Jordan chain, and this means that the matrix A had a rather special form. We'll give a name to matrices of this type:

Definition 2.6.5. We define the *Jordan block* of degree k with eigenvalue λ to be the $k \times k$ matrix $J_{\lambda,k}$ whose entries are given by

$$\gamma_{ij} = \begin{cases} \lambda & \text{if } j = i \\ 1 & \text{if } j = i + 1 \\ 0 & \text{otherwise.} \end{cases}$$

So, for example,

$$J_{1,2} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}, \quad J_{4i-7,3} = \begin{pmatrix} 4i-7 & 1 & 0 \\ 0 & 4i-7 & 1 \\ 0 & 0 & 4i-7 \end{pmatrix}, \quad \text{and} \quad J_{0,4} = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

2 The Jordan Canonical Form

are Jordan blocks.

It should be clear that the matrix of T with respect to the basis $\mathbf{v}_1, \dots, \mathbf{v}_n$ of K^n is a Jordan block of degree n if and only if $\mathbf{v}_1, \dots, \mathbf{v}_n$ is a Jordan chain for A .

Note that the minimal polynomial of $J_{\lambda,k}$ is equal to $(x - \lambda)^k$, and the characteristic polynomial is $(\lambda - x)^k$.

Warning. Some authors put the 1's below rather than above the main diagonal in a Jordan block. This corresponds to writing the Jordan chain in reverse order. This is an arbitrary choice but in this course we stick to our convention - when you read other notes/books be careful to check which convention they use.

2.7 Jordan bases and the Jordan canonical form

Definition 2.7.1. A *Jordan basis* for A is a basis of K^n consisting of one or more Jordan chains strung together.

Such a basis will look like

$$w_{11}, \dots, w_{1k_1}, w_{21}, \dots, w_{2k_2}, \dots, w_{s1}, \dots, w_{sk_s},$$

where, for $1 \leq i \leq s$, w_{i1}, \dots, w_{ik_i} is a Jordan chain (for some eigenvalue λ_i).

This definition is the key to defining Jordan canonical form. However, until we prove the main theorem (Theorem 2.1) we do not know that such bases always exist! At least Lemma 2.6.3 shows us that the matrices representing the linear map of a single Jordan block has such a basis, which gives us hope that this definition will be fruitful.

We denote the $m \times n$ matrix in which all entries are 0 by $\mathbf{0}_{m,n}$. If A is an $m \times m$ matrix and B an $n \times n$ matrix, then we denote the $(m + n) \times (m + n)$ matrix with block form

$$\left(\begin{array}{c|c} A & \mathbf{0}_{m,n} \\ \hline \mathbf{0}_{n,m} & B \end{array} \right),$$

by $A \oplus B$, the *direct sum* of A and B . For example

$$\begin{pmatrix} -1 & 2 \\ 0 & 1 \end{pmatrix} \oplus \begin{pmatrix} 1 & 1 & -1 \\ 1 & 0 & 1 \\ 2 & 0 & -2 \end{pmatrix} = \begin{pmatrix} -1 & 2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & -1 \\ 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 2 & 0 & -2 \end{pmatrix}.$$

It's clear that the matrix of T with respect to a Jordan basis is the direct sum $J_{\lambda_1, k_1} \oplus J_{\lambda_2, k_2} \oplus \dots \oplus J_{\lambda_s, k_s}$ of the corresponding Jordan blocks.

You'll be asked to prove the following lemma on an examples sheet, which follows from the definition of the direct sum and the characteristic and minimal polynomials.

2 The Jordan Canonical Form

Lemma 2.7.2. *Suppose that $M = A \oplus B$. Then the characteristic polynomial $c_M(x)$ is the product of $c_A(x)$ and $c_B(x)$, and the minimal polynomial $\mu_M(x)$ is the lowest common multiple of $\mu_A(x)$ and $\mu_B(x)$.*

It is now time for us to state the main theorem of this section, which says that if K is the complex numbers \mathbb{C} , then Jordan bases exist.

Theorem 2.7.3. *Let A be an $n \times n$ matrix over \mathbb{C} . Then there exists a Jordan basis for A , and hence A is similar to a matrix J which is a direct sum of Jordan blocks. The Jordan blocks occurring in J are uniquely determined by A .*

The matrix J in the theorem is said to be the *Jordan canonical form (JCF)* or sometimes *Jordan normal form* of A . It is uniquely determined by A up to the order of the blocks.

Remark. *The only reason we need $K = \mathbb{C}$ in this theorem is to ensure that A has at least one eigenvalue. If $K = \mathbb{R}$ (or \mathbb{Q}), we'd run into trouble with $\begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix}$; this matrix has no eigenvalues, since $c_A(x) = x^2 + 1$ has no roots in K . So it certainly has no Jordan chains. The theorem is valid more generally for any field K which is such that any non-constant polynomial in $K[x]$ has a root in K (one calls such fields algebraically closed; there are many more of them out there than just \mathbb{C}).*

We will prove the theorem later, in Section 2.11. First we derive some consequences and study methods for calculating the JCF of a matrix. Note that, by Theorem 1.5.2, if P is the matrix having the Jordan basis as columns, then $P^{-1}AP = J$.

Theorem 2.7.4 (Consequences of the JCF). *Let $A \in \mathbb{C}^{n,n}$, and $\{\lambda_1, \dots, \lambda_r\}$ be the set of eigenvalues of A .*

(i) *The characteristic polynomial of A is*

$$(-1)^n \prod_{i=1}^r (x - \lambda_i)^{a_i},$$

where a_i is the sum of the degrees of the Jordan blocks of A of eigenvalue λ_i .

(ii) *The minimal polynomial of A is*

$$\prod_{i=1}^r (x - \lambda_i)^{b_i},$$

where b_i is the largest among the degrees of the Jordan blocks of A of eigenvalue λ_i .

(iii) *A is diagonalizable if and only if $\mu_A(x)$ has no repeated factors.*

Proof. We know that the characteristic and minimal polynomials of A and J , its JCF, are the same. So the first two parts follow from applying Lemma 2.7.2 (multiple times) to J . For the last part, notice that if A is diagonalizable, the JCF of A is just the diagonal form of A ; since the JCF is unique, it follows that A is diagonalizable if and only if every Jordan block for A has size 1, so all of the numbers b_i are 1. □

2.8 The JCF when $n=2$ and 3

When $n = 2$ and $n = 3$, the JCF can be deduced just from the minimal and characteristic polynomials. Let us consider these cases.

When $n = 2$, we have either two distinct eigenvalues λ_1, λ_2 , or a single repeated eigenvalue λ_1 . If the eigenvalues are distinct, then by Corollary 2.2.3 A is diagonalizable and the JCF is the diagonal matrix $J_{\lambda_1,1} \oplus J_{\lambda_2,1}$.

Example 3. $A = \begin{pmatrix} 1 & 4 \\ 1 & 1 \end{pmatrix}$. We calculate $c_A(x) = x^2 - 2x - 3 = (x - 3)(x + 1)$, so there are two distinct eigenvalues, 3 and -1 . Associated eigenvectors are $\begin{pmatrix} 2 \\ 1 \end{pmatrix}$ and $\begin{pmatrix} -2 \\ 1 \end{pmatrix}$, so we put $P = \begin{pmatrix} 2 & -2 \\ 1 & 1 \end{pmatrix}$ and then $P^{-1}AP = \begin{pmatrix} 3 & 0 \\ 0 & -1 \end{pmatrix}$.

If the eigenvalues are equal, then there are two possible JCFs, $J_{\lambda_1,1} \oplus J_{\lambda_1,1}$, which is a scalar matrix, and $J_{\lambda_1,2}$. The minimal polynomial is respectively $(x - \lambda_1)$ and $(x - \lambda_1)^2$ in these two cases. In fact, these cases can be distinguished without any calculation whatsoever, because in the first case A is a scalar multiple of the identity, and in particular A is already in JCF.

In the second case, a Jordan basis consists of a single Jordan chain of length 2. To find such a chain, let \mathbf{v}_2 be any vector for which $(A - \lambda_1 I_2)\mathbf{v}_2 \neq \mathbf{0}$ and let $\mathbf{v}_1 = (A - \lambda_1 I_2)\mathbf{v}_2$. (Note that, in practice, it is often easier to find the vectors in a Jordan chain in reverse order.)

Example 4. $A = \begin{pmatrix} 1 & 4 \\ -1 & -3 \end{pmatrix}$. We have $c_A(x) = x^2 + 2x + 1 = (x + 1)^2$, so there is a single eigenvalue -1 with multiplicity 2. Since the first column of $A + I_2$ is non-zero, we can choose $\mathbf{v}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ and $\mathbf{v}_1 = (A + I_2)\mathbf{v}_2 = \begin{pmatrix} 2 \\ -1 \end{pmatrix}$, so $P = \begin{pmatrix} 2 & 1 \\ -1 & 0 \end{pmatrix}$ and $P^{-1}AP = \begin{pmatrix} -1 & 1 \\ 0 & -1 \end{pmatrix}$.

Now let $n = 3$. If there are three distinct eigenvalues, then A is diagonalizable.

Suppose that there are two distinct eigenvalues, so one has multiplicity 2, and the other has multiplicity 1. Let the eigenvalues be $\lambda_1, \lambda_1, \lambda_2$, with $\lambda_1 \neq \lambda_2$. Then there are two possible JCFs for A , $J_{\lambda_1,1} \oplus J_{\lambda_1,1} \oplus J_{\lambda_2,1}$ and $J_{\lambda_1,2} \oplus J_{\lambda_2,1}$, and the minimal polynomial is $(x - \lambda_1)(x - \lambda_2)$ in the first case and $(x - \lambda_1)^2(x - \lambda_2)$ in the second.

In the first case, a Jordan basis is a union of three Jordan chains of length 1, each of which consists of an eigenvector of A .

Example 5. $A = \begin{pmatrix} 2 & 0 & 0 \\ 1 & 5 & 2 \\ -2 & -6 & -2 \end{pmatrix}$. Then

$$c_A(x) = (2 - x)[(5 - x)(-2 - x) + 12] = (2 - x)(x^2 - 3x + 2) = (2 - x)^2(1 - x).$$

We know from the theory above that the minimal polynomial must be $(x - 2)(x - 1)$ or $(x - 2)^2(x - 1)$. We can decide which simply by calculating $(A - 2I_3)(A - I_3)$ to test whether or not

2 The Jordan Canonical Form

it is 0. We have

$$A - 2I_3 = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 3 & 2 \\ -2 & -6 & -4 \end{pmatrix}, \quad A - I_3 = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 4 & 2 \\ -2 & -6 & -3 \end{pmatrix},$$

and the product of these two matrices is 0, so $\mu_A = (x - 2)(x - 1)$.

The eigenvectors \mathbf{v} for $\lambda_1 = 2$ satisfy $(A - 2I_3)\mathbf{v} = \mathbf{0}$, and we must find two linearly independent solutions; for example we can take $\mathbf{v}_1 = \begin{pmatrix} 0 \\ 2 \\ -3 \end{pmatrix}$, $\mathbf{v}_2 = \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}$. An eigenvector for the

eigenvalue 1 is $\mathbf{v}_3 = \begin{pmatrix} 0 \\ 1 \\ -2 \end{pmatrix}$, so we can choose

$$P = \begin{pmatrix} 0 & 1 & 0 \\ 2 & -1 & 1 \\ -3 & 1 & -2 \end{pmatrix}$$

and then $P^{-1}AP$ is diagonal with entries 2, 2, 1.

In the second case, there are two Jordan chains, one for λ_1 of length 2, and one for λ_2 of length 1. For the first chain, we need to find a vector \mathbf{v}_2 with $(A - \lambda_1 I_3)^2 \mathbf{v}_2 = \mathbf{0}$ but $(A - \lambda_1 I_3) \mathbf{v}_2 \neq \mathbf{0}$, and then the chain is $\mathbf{v}_1 = (A - \lambda_1 I_3) \mathbf{v}_2, \mathbf{v}_2$. For the second chain, we simply need an eigenvector for λ_2 .

Example 6. $A = \begin{pmatrix} 3 & 2 & 1 \\ 0 & 3 & 1 \\ -1 & -4 & -1 \end{pmatrix}$. Then

$$c_A(x) = (3 - x)[(3 - x)(-1 - x) + 4] - 2 + (3 - x) = -x^3 + 5x^2 - 8x + 4 = (2 - x)^2(1 - x),$$

as in Example 3. We have

$$A - 2I_3 = \begin{pmatrix} 1 & 2 & 1 \\ 0 & 1 & 1 \\ -1 & -4 & -3 \end{pmatrix}, \quad (A - 2I_3)^2 = \begin{pmatrix} 0 & 0 & 0 \\ -1 & -3 & -2 \\ 2 & 6 & 4 \end{pmatrix}, \quad (A - I_3) = \begin{pmatrix} 2 & 2 & 1 \\ 0 & 2 & 1 \\ -1 & -4 & -2 \end{pmatrix}.$$

and we can check that $(A - 2I_3)(A - I_3)$ is non-zero, so we must have $\mu_A = (x - 2)^2(x - 1)$.

For the Jordan chain of length 2, we need a vector with $(A - 2I_3)^2 \mathbf{v}_2 = \mathbf{0}$ but $(A - 2I_3) \mathbf{v}_2 \neq \mathbf{0}$, and we can choose $\mathbf{v}_2 = \begin{pmatrix} 2 \\ 0 \\ -1 \end{pmatrix}$. Then $\mathbf{v}_1 = (A - 2I_3) \mathbf{v}_2 = \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}$. An eigenvector for the

eigenvalue 1 is $\mathbf{v}_3 = \begin{pmatrix} 0 \\ 1 \\ -2 \end{pmatrix}$, so we can choose

$$P = \begin{pmatrix} 1 & 2 & 0 \\ -1 & 0 & 1 \\ 1 & -1 & -2 \end{pmatrix}$$

2 The Jordan Canonical Form

and then

$$P^{-1}AP = \begin{pmatrix} 2 & 1 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

Finally, suppose that there is a single eigenvalue, λ_1 , so $c_A = (\lambda_1 - x)^3$. There are three possible JCFs for A , namely $J_{\lambda_1,1} \oplus J_{\lambda_1,1} \oplus J_{\lambda_1,1}$, $J_{\lambda_1,2} \oplus J_{\lambda_1,1}$, and $J_{\lambda_1,3}$, and the minimal polynomials in the three cases are $(x - \lambda_1)$, $(x - \lambda_1)^2$, and $(x - \lambda_1)^3$, respectively.

In the first case, J is a scalar matrix, and $A = PJP^{-1} = J$, so this is recognisable immediately.

In the second case, there are two Jordan chains, one of length 2 and one of length 1. For the first, we choose \mathbf{v}_2 with $(A - \lambda_1 I_3)\mathbf{v}_2 \neq \mathbf{0}$, and let $\mathbf{v}_1 = (A - \lambda_1 I_3)\mathbf{v}_2$. (This case is easier than the case illustrated in Example 4, because we have $(A - \lambda_1 I_3)^2 \mathbf{v} = \mathbf{0}$ for all $\mathbf{v} \in \mathbb{C}^{3,1}$.) For the second Jordan chain, we choose \mathbf{v}_3 to be an eigenvector for λ_1 such that \mathbf{v}_1 and \mathbf{v}_3 are linearly independent.

Example 7. $A = \begin{pmatrix} 0 & 2 & 1 \\ -1 & -3 & -1 \\ 1 & 2 & 0 \end{pmatrix}$. Then

$$c_A(x) = -x[(3+x)x+2] - 2(x+1) - 2 + (3+x) = -x^3 - 3x^2 - 3x - 1 = -(1+x)^3.$$

We have

$$A + I_3 = \begin{pmatrix} 1 & 2 & 1 \\ -1 & -2 & -1 \\ 1 & 2 & 1 \end{pmatrix},$$

and we can check that $(A + I_3)^2 = \mathbf{0}$. The first column of $A + I_3$ is non-zero, so $(A + I_3) \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \neq \mathbf{0}$,

and we can choose $\mathbf{v}_2 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$ and $\mathbf{v}_1 = (A + I_3)\mathbf{v}_2 = \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}$. For \mathbf{v}_3 we need to choose a

vector which is not a multiple of \mathbf{v}_1 such that $(A + I_3)\mathbf{v}_3 = \mathbf{0}$, and we can choose $\mathbf{v}_3 = \begin{pmatrix} 0 \\ 1 \\ -2 \end{pmatrix}$.

So we have

$$P = \begin{pmatrix} 1 & 1 & 0 \\ -1 & 0 & 1 \\ 1 & 0 & -2 \end{pmatrix}$$

and then

$$P^{-1}AP = \begin{pmatrix} -1 & 1 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{pmatrix}.$$

In the third case, there is a single Jordan chain, and we choose \mathbf{v}_3 such that $(A - \lambda_1 I_3)^2 \mathbf{v}_3 \neq \mathbf{0}$, $\mathbf{v}_2 = (A - \lambda_1 I_3)\mathbf{v}_3$, $\mathbf{v}_1 = (A - \lambda_1 I_3)^2 \mathbf{v}_3$.

2 The Jordan Canonical Form

Example 8. $A = \begin{pmatrix} 0 & 1 & 0 \\ -1 & -1 & 1 \\ 1 & 0 & -2 \end{pmatrix}$. Then

$$c_A(x) = -x[(2+x)(1+x)] - (2+x) + 1 = -(1+x)^3.$$

We have

$$A + I_3 = \begin{pmatrix} 1 & 1 & 0 \\ -1 & 0 & 1 \\ 1 & 0 & -1 \end{pmatrix}, \quad (A + I_3)^2 = \begin{pmatrix} 0 & 1 & 1 \\ 0 & -1 & -1 \\ 0 & 1 & 1 \end{pmatrix},$$

so $(A + I_3)^2 \neq 0$ and $\mu_A = (x + 1)^3$. For \mathbf{v}_3 , we need a vector that is not in the nullspace of $(A + I_3)^2$. Since the second column, which is the image of $\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$ is non-zero, we can choose

$\mathbf{v}_3 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$, and then $\mathbf{v}_2 = (A + I_3)\mathbf{v}_3 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$ and $\mathbf{v}_1 = (A + I_3)\mathbf{v}_2 = \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}$. So we have

$$P = \begin{pmatrix} 1 & 1 & 0 \\ -1 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix}$$

and then

$$P^{-1}AP = \begin{pmatrix} -1 & 1 & 0 \\ 0 & -1 & 1 \\ 0 & 0 & -1 \end{pmatrix}.$$

2.9 The general case

In the examples above, we could tell what the sizes of the Jordan blocks were for each eigenvalue from the dimensions of the eigenspaces, since the dimension of the eigenspace for each eigenvalue λ is the number of blocks for that eigenvalue. This doesn't work for $n = 4$: for instance, the matrices

$$A_1 = J_{\lambda,2} \oplus J_{\lambda,2}$$

and

$$A_2 = J_{\lambda,3} \oplus J_{\lambda,1}$$

both have only one eigenvalue (λ) with the eigenspace being of dimension 2.

(Knowing the minimal polynomial helps, but it's a bit of a pain to calculate – generally the easiest way to find the minimal polynomial is to calculate the JCF first! Worse still, it still doesn't uniquely determine the JCF in large dimensions, since

$$A_3 = J_{\lambda,3} \oplus J_{\lambda,3} \oplus J_{\lambda,1}$$

and

$$A_4 = J_{\lambda,3} \oplus J_{\lambda,2} \oplus J_{\lambda,2}$$

2 The Jordan Canonical Form

have the same minimal polynomial, the same characteristic polynomial, and the same number of blocks.)

In general, we can compute the JCF from the dimensions of the generalised eigenspaces. Notice that the matrices A_1 and A_2 can be distinguished by looking at the dimensions of their generalised eigenspaces: the generalised eigenspace for λ of index 2 has dimension 4 for A_1 (it's the whole space) but dimension only 3 for A_2 .

Theorem 2.9.1. *Let λ be an eigenvalue of a matrix $A \in \mathbb{C}^{n,n}$, and let J be the JCF of A . Then*

- (i) *The number of Jordan blocks of J with eigenvalue λ is equal to $\text{nullity}(A - \lambda I_n)$.*
- (ii) *More generally, for $i > 0$, the number of Jordan blocks of J with eigenvalue λ and degree at least i is equal to $\text{nullity}((A - \lambda I_n)^i) - \text{nullity}((A - \lambda I_n)^{i-1})$.*

Note that this proves the uniqueness part of Theorem 2.7.3: the theorem says that the block sizes of the Jordan form of A are determined by the dimensions of the generalised eigenspaces of A for each eigenvalue, so any two Jordan canonical forms for A must have the same blocks (possibly ordered differently).

Proof. By Proposition 2.6.4, the corresponding generalised eigenspaces of A and J have the same dimensions, so we may assume WLOG that $A = J$. So A is a direct sum of several Jordan blocks $J_{\lambda_1, k_1} \oplus \cdots \oplus J_{\lambda_s, k_s}$.

However, it's easy to see that the dimension of the generalised λ -eigenspace of index i of a direct sum $A \oplus B$ is the sum of the dimensions of the generalised λ eigenspaces of index i of A and of B . Hence it suffices to prove the theorem for a single Jordan block $J_{\lambda, k}$.

But we know that $(J_{\lambda, k} - \lambda I_k)^i$ has a single diagonal line of ones i places above the diagonal, for $i < k$, and is 0 for $i \geq k$. Hence the dimension of its kernel is i for $0 \leq i \leq k$ and k for $i \geq k$. This clearly implies the theorem when A is a single Jordan block, and hence for any A . \square

2.10 Examples

Example 9. $A = \begin{pmatrix} -1 & -3 & -1 & 0 \\ 0 & 2 & 1 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 3 & 1 & -1 \end{pmatrix}$. Then $c_A(x) = (-1 - x)^2(2 - x)^2$, so there are two

eigenvalues $-1, 2$, both with multiplicity 2. There are four possibilities for the JCF (one or two blocks for each of the two eigenvalues). We could determine the JCF by computing the minimal polynomial μ_A but it is probably easier to compute the nullities of the eigenspaces and use

2 The Jordan Canonical Form

Theorem 2.9.1. We have

$$A + I_4 = \begin{pmatrix} 0 & -3 & -1 & 0 \\ 0 & 3 & 1 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 3 & 1 & 0 \end{pmatrix}, \quad (A - 2I_4) = \begin{pmatrix} -3 & -3 & -1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 3 & 1 & -3 \end{pmatrix},$$

$$(A - 2I_4)^2 = \begin{pmatrix} 9 & 9 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & -9 & 0 & 9 \end{pmatrix}.$$

The rank of $A + I_4$ is clearly 2, so its nullity is also 2, and hence there are two Jordan blocks with eigenvalue -1 . The three non-zero rows of $(A - 2I_4)$ are linearly independent, so its rank is 3, hence its nullity 1, so there is just one Jordan block with eigenvalue 2, and the JCF of A is $J_{-1,1} \oplus J_{-1,1} \oplus J_{2,2}$.

For the two Jordan chains of length 1 for eigenvalue -1 , we just need two linearly independent

eigenvectors, and the obvious choice is $\mathbf{v}_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$, $\mathbf{v}_2 = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$. For the Jordan chain $\mathbf{v}_3, \mathbf{v}_4$ for

eigenvalue 2, we need to choose \mathbf{v}_4 in the nullspace of $(A - 2I_4)^2$ but not in the nullspace of $A - 2I_4$. (This is why we calculated $(A - 2I_4)^2$.) An obvious choice here is $\mathbf{v}_4 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}$, and then

$\mathbf{v}_3 = \begin{pmatrix} -1 \\ 1 \\ 0 \\ 1 \end{pmatrix}$, and to transform A to JCF, we put

$$P = \begin{pmatrix} 1 & 0 & -1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix}, \quad P^{-1} = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & -1 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}, \quad P^{-1}AP = \begin{pmatrix} -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 2 & 1 \\ 0 & 0 & 0 & 2 \end{pmatrix}.$$

Example 10. $A = \begin{pmatrix} -2 & 0 & 0 & 0 \\ 0 & -2 & 1 & 0 \\ 0 & 0 & -2 & 0 \\ 1 & 0 & -2 & -2 \end{pmatrix}$. Then $c_A(x) = (-2 - x)^4$, so there is a single eigen-

value -2 with multiplicity 4. We find $(A + 2I_4) = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & 0 & -2 & 0 \end{pmatrix}$, and $(A + 2I_4)^2 = 0$, so

$\mu_A = (x + 2)^2$, and the JCF of A could be $J_{-2,2} \oplus J_{-2,2}$ or $J_{-2,2} \oplus J_{-2,1} \oplus J_{-2,1}$.

2 The Jordan Canonical Form

To decide which case holds, we calculate the nullity of $A + 2I_4$ which, by Theorem 2.9.1, is equal to the number of Jordan blocks with eigenvalue -2 . Since $A + 2I_4$ has just two non-zero rows, which are distinct, its rank is clearly 2, so its nullity is $4 - 2 = 2$, and hence the JCF of A is $J_{-2,2} \oplus J_{-2,2}$.

A Jordan basis consists of a union of two Jordan chains, which we will call $\mathbf{v}_1, \mathbf{v}_2$, and $\mathbf{v}_3, \mathbf{v}_4$, where \mathbf{v}_1 and \mathbf{v}_3 are eigenvectors and \mathbf{v}_2 and \mathbf{v}_4 are generalised eigenvectors of index 2. To find such chains, it is probably easiest to find \mathbf{v}_2 and \mathbf{v}_4 first and then to calculate $\mathbf{v}_1 = (A + 2I_4)\mathbf{v}_2$ and $\mathbf{v}_3 = (A + 2I_4)\mathbf{v}_4$.

Although it is not hard to find \mathbf{v}_2 and \mathbf{v}_4 in practice, we have to be careful, because they need to be chosen so that no linear combination of them lies in the nullspace of $(A + 2I_4)$. In fact, since this nullspace is spanned by the second and fourth standard basis vectors, the obvious choice is

$$\mathbf{v}_2 = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \mathbf{v}_4 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix}, \text{ and then } \mathbf{v}_1 = (A + 2I_4)\mathbf{v}_2 = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}, \mathbf{v}_3 = (A + 2I_4)\mathbf{v}_4 = \begin{pmatrix} 0 \\ 1 \\ 0 \\ -2 \end{pmatrix}, \text{ so to}$$

transform A to JCF, we put

$$P = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & -2 & 0 \end{pmatrix}, \quad P^{-1} = \begin{pmatrix} 0 & 2 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}, \quad P^{-1}AP = \begin{pmatrix} -2 & 1 & 0 & 0 \\ 0 & -2 & 0 & 0 \\ 0 & 0 & -2 & 1 \\ 0 & 0 & 0 & -2 \end{pmatrix}.$$

2.11 Proof of Theorem 2.7.3

We proceed by induction on $n = \dim(V)$. The case $n = 1$ is clear.

We are looking for a vector space of dimension less than n , related to T to apply our inductive hypothesis to. Let λ be an eigenvalue of T and set $S := T - \lambda I_V$. Then we let $U = \text{im}(S)$ and $m = \dim(U)$. Using the Rank-Nullity Theorem we see that $m = \text{rank}(S) = n - \text{nullity}(S) < n$, because there exists at least one eigenvector of T for λ , which lies in the nullspace of $S = T - \lambda I_V$. For $\mathbf{u} \in U$, we have $\mathbf{u} = S(\mathbf{v})$ for some $\mathbf{v} \in V$, and hence $T(\mathbf{u}) = TS(\mathbf{v}) = ST(\mathbf{v}) \in \text{im}(S) = U$. Note that $TS = ST$ because $T(T - \lambda I_V) = T^2 - T\lambda I_V = T^2 - \lambda I_V T = (T - \lambda I_V)T$. So T maps U to U and thus T restricts to a linear map $T_U : U \rightarrow U$. Since $m < n$, we can apply our inductive hypothesis to T_U to deduce that U has a basis $\mathbf{e}_1, \dots, \mathbf{e}_m$, which is a disjoint union of Jordan chains for T_U (for all eigenvalues of T_U).

It is our job to show how to extend this Jordan basis of U to one of V . We do this in two stages. Firstly, let $\mathbf{v}_1, \dots, \mathbf{v}_k$ be one of the l disjoint Jordan chains for eigenvalue λ for T_U (where l could be 0), so we have $T(\mathbf{v}_1) = T_U(\mathbf{v}_1) = \lambda\mathbf{v}_1$, $T(\mathbf{v}_i) = T_U(\mathbf{v}_i) = \lambda\mathbf{v}_i + \mathbf{v}_{i-1}$, $2 \leq i \leq k$. Now, since $\mathbf{v}_k \in U = \text{im } S = \text{im}(T - \lambda I_V)$, we can find $\mathbf{v}_{k+1} \in V$ with $T(\mathbf{v}_{k+1}) = \lambda\mathbf{v}_{k+1} + \mathbf{v}_k$, thereby extending the chain by an extra vector of V .

We do this for each of the l disjoint chains for λ and so at this point we have adjoined l new vectors to the basis. Let us call these new vectors $\mathbf{w}_1, \dots, \mathbf{w}_l$.

For the second stage, observe that the first vector in each of the l chains lies in the eigenspace

2 The Jordan Canonical Form

of T_U for λ . We know that the dimension of the eigenspace of T for λ is the dimension of the nullspace of S , which is $n - m$. So we can adjoin $(n - m) - l$ (which could be 0) further eigenvectors of T to the l that we have already to complete a basis of the nullspace of $(T - \lambda I_V)$. Let us call these $(n - m) - l$ new vectors $\mathbf{w}_{l+1}, \dots, \mathbf{w}_{n-m}$. They are adjoined to our basis of V in the second stage. They each form a Jordan chain of length 1 (since they are not in the image of $S = T - \lambda I_V$), so we now have a collection of n vectors which form a disjoint union of Jordan chains.

To complete the proof, we need to show that these n vectors form a basis of V , for which it is enough to show that they are linearly independent.

Suppose that

$$\alpha_1 \mathbf{w}_1 + \dots + \alpha_{n-m} \mathbf{w}_{n-m} + \mathbf{x} = \mathbf{0}, \quad \text{where } \mathbf{x} = \beta_1 \mathbf{e}_1 + \dots + \beta_m \mathbf{e}_m, \quad (3)$$

a linear combination of the basis vectors $\mathbf{e}_1, \dots, \mathbf{e}_m$ of U . We now apply S to both sides of this equation, recalling that $S(\mathbf{w}_{l+i}) = \mathbf{0}$ for $i \geq 1$, by definition.

$$\alpha_1 S(\mathbf{w}_1) + \dots + \alpha_l S(\mathbf{w}_l) + S(\mathbf{x}) = \mathbf{0}. \quad (4)$$

By the construction of the \mathbf{w}_i , each of the $S(\mathbf{w}_i)$ for $1 \leq i \leq l$ is the last member of one of the l Jordan chains for T_U . Let this set of l vectors \mathbf{e}_j be $L = \{j \mid \mathbf{e}_j = S(\mathbf{w}_i) \text{ for some } 1 \leq i \leq l\}$. Now examine the last term

$$S(\mathbf{x}) = (T - \lambda I_n)(\mathbf{x}) = (T_U - \lambda I_m)(\mathbf{x}) = \beta_1 (T_U - \lambda I_m)(\mathbf{e}_1) + \dots + \beta_m (T_U - \lambda I_m)(\mathbf{e}_m).$$

Each $(T_U - \lambda I_m)(\mathbf{e}_j)$ is a linear combination of the basis vectors of U from the subset

$$\{\mathbf{e}_1, \dots, \mathbf{e}_m\} \setminus \{\mathbf{e}_j \mid j \in L\}.$$

Indeed, this follows because after application of S we must have ‘moved’ down our Jordan chains for T_U . It now follows from the linear independence of the basis $\mathbf{e}_1, \dots, \mathbf{e}_m$, that $\alpha_i = 0$ for all $1 \leq i \leq l$.

So Equation (4) is now just

$$S(\mathbf{x}) = \mathbf{0},$$

and so \mathbf{x} is in the eigenspace of T_U for the eigenvalue λ . Equation (3) looks like

$$\alpha_{l+1} \mathbf{w}_{l+1} + \dots + \alpha_{n-m} \mathbf{w}_{n-m} + \mathbf{x} = \mathbf{0}. \quad (5)$$

By construction, $\mathbf{w}_{l+1}, \dots, \mathbf{w}_{n-m}$ extend a basis of the eigenspace of T_U to a basis of the eigenspace of T for λ . Lemma 1.2.2 now applies (to the eigenspace of T), yielding $\alpha_i = 0$ for $l + 1 \leq i \leq n - m$ and $\mathbf{x} = \mathbf{0}$. Since $\mathbf{e}_1, \dots, \mathbf{e}_m$ is a basis for U , we must have all $\beta_j = 0$, which completes the proof.

2 The Jordan Canonical Form

2.12 An algorithm to compute the Jordan canonical form in general (brute force)

Whereas the examples above in Sections 2.8 and 2.10 explain some shortcuts, tricks and computational recipes to compute, given a matrix $A \in \mathbb{C}^{n,n}$, a Jordan canonical form J for A as well as a matrix P (invertible) such that $J = P^{-1}AP$, it may also be useful to know how this can be done systematically, *provided* we know all the eigenvalues, $\lambda_1, \dots, \lambda_s$, say, of A .

Algorithm:

Step 1: Compute J . This amounts to knowing, for a given eigenvalue λ , the number of Jordan blocks of degree/size i in J . By Theorem 2.9.1, (ii), this number is

$$\begin{aligned} & (\dim N_i(A, \lambda) - \dim N_{i-1}(A, \lambda)) - (\dim N_{i+1}(A, \lambda) - \dim N_i(A, \lambda)) \\ & = 2 \dim N_i(A, \lambda) - \dim N_{i-1}(A, \lambda) - \dim N_{i+1}(A, \lambda). \end{aligned}$$

So the computation of J is no problem then.

Step 2: Compute P . You can proceed as follows: pick an eigenvalue λ . Now suppose

$$N_1 \geq N_2 \geq \dots \geq N_r$$

are the sizes of the Jordan blocks with eigenvalue λ (repeats among the N_i allowed if there are several blocks of the same size; we order them according to decreasing size for definiteness). Then pick a vector $v_{1,1} \in V$ with

$$(A - \lambda I_n)^{N_1} v_{1,1} = 0, (A - \lambda I_n)^{N_1-1} v_{1,1} \neq 0$$

(note that this amounts to solving several systems of linear equations ultimately- we leave the details of how to accomplish this step to you). Then put

$$v_{1,2} := (A - \lambda I_n)v_{1,1}, v_{1,3} := (A - \lambda I_n)^2 v_{1,1}, \dots, v_{1,N_1} := (A - \lambda I_n)^{N_1-1} v_{1,1}.$$

Note that $(v_{1,N_1}, \dots, v_{1,1})$ is then a Jordan chain. If $r = 1$, we are done, else we choose a vector $v_{2,1} \in V$ with

$$(A - \lambda I_n)^{N_2} v_{2,1} = 0, (A - \lambda I_n)^{N_2-1} v_{2,1} \notin \langle v_{1,1}, \dots, v_{1,N_1} \rangle.$$

So note that the second condition has become more restrictive: we want that $(A - \lambda I_n)^{N_2-1} v_{2,1}$ is not just nonzero, but not in the span $V_1 := \langle v_{1,1}, \dots, v_{1,N_1} \rangle$ of the first bunch of basis vectors. Equivalently, we want it to be nonzero in the quotient V/V_1 , for those of you who know what quotient vector spaces are (which isn't required). We then put

$$v_{2,2} := (A - \lambda I_n)v_{2,1}, v_{2,3} := (A - \lambda I_n)^2 v_{2,1}, \dots, v_{2,N_2} := (A - \lambda I_n)^{N_2-1} v_{2,1}.$$

Then by construction $(v_{2,N_2}, \dots, v_{2,1})$ is a Jordan chain, and $v_{1,1}, \dots, v_{1,N_1}, v_{2,1}, \dots, v_{2,N_2}$ are linearly independent (for those who know quotient spaces, an easy way to check this is to notice that $v_{2,N_2}, \dots, v_{2,1}$ are a Jordan chain in V/V_1). If $r = 2$, we are done, otherwise we continue in the same fashion: pick $v_{3,1} \in V$ with

$$(A - \lambda I_n)^{N_3} v_{3,1} = 0, (A - \lambda I_n)^{N_3-1} v_{3,1} \notin \langle v_{1,1}, \dots, v_{1,N_1}, v_{2,1}, \dots, v_{2,N_2} \rangle,$$

2 The Jordan Canonical Form

and now you should see what the pattern to continue is. Finally you end up with vectors

$$v_{1,1}, \dots, v_{1,N_1}, v_{2,1}, \dots, v_{2,N_2}, \dots, v_{r,1}, \dots, v_{r,N_r} \in \mathbb{C}^n.$$

Listing these in reverse order gives us the first $N_1 + \dots + N_r$ columns of P . Now we repeat the same procedure for the remaining eigenvalues of A other than λ , adding a bunch of columns to P at each step in this way. That gives us the desired base change matrix P .

2.13 Grand finale

At this point we would like to take a step back and formulate the **basic facts of the spectral theory of matrices** we have obtained so far in a way that is both easy to remember and convenient to use in many applications. We use the more standard $\mathbb{C}^{n \times n}$ for $\mathbb{C}^{n,n}$ and \mathbb{C}^n for $\mathbb{C}^{n,1}$ below.

Theorem 2.13.1. *Let $A \in \mathbb{C}^{n \times n}$ be a square matrix with complex entries, $p \in \mathbb{C}[x]$ any polynomial. Then if λ is an eigenvalue of A , $p(\lambda)$ is an eigenvalue of $p(A)$, and any eigenvalue of $p(A)$ is of this form.*

In fact, you have shown that as an exercise on the first assignment.

Theorem 2.13.2. *For $A \in \mathbb{C}^{n \times n}$ let*

$$N_i(A, \lambda)$$

be the null-space of $(A - \lambda I_n)^i$, so non-zero elements in $N_i(A, \lambda)$ are generalised eigenvectors of A w.r.t. λ of index i (index 1 being genuine eigenvectors). Then every vector in \mathbb{C}^n can be written as a sum of eigenvectors of A , genuine or generalised.

This follows immediately from Theorem [2.7.3](#).

Theorem 2.13.3. (i) *Suppose $A, B \in \mathbb{C}^{n \times n}$ are similar, in the sense that there exists an invertible $n \times n$ matrix S with $B = S^{-1}AS$. Then A and B have the same set of eigenvalues:*

$$\lambda_1 = \mu_1, \dots, \lambda_k = \mu_k$$

(here the λ 's are the eigenvalues for A , the μ 's the ones for B), and in addition we have

$$(*) \quad \dim N_i(A, \lambda_j) = \dim N_i(B, \mu_j)$$

for all i, j .

(ii) *Conversely, if $A, B \in \mathbb{C}^{n \times n}$ have the same eigenvalues $\lambda_1 = \mu_1, \dots, \lambda_k = \mu_k$ as above, and $(*)$ holds for all i and j , then A and B are similar.*

2 The Jordan Canonical Form

Whereas (i) is obvious (but also the content of Proposition 2.6.4), (ii) follows from the uniqueness part of Theorem 2.7.3 (or Theorem 2.9.1 alternatively).

The three results above are the basic results of spectral theory, in some sense even more basic than the Jordan canonical form itself. Also clearly

$$N_1(A, \lambda) \subset N_2(A, \lambda) \subset N_3(A, \lambda) \subset \dots$$

and denoting by $d(\lambda)$ the smallest index from which these spaces are equal to each other (the index of the eigenvalue λ), we have: if $\lambda_1, \dots, \lambda_k$ are the distinct eigenvalues of A , we have for the minimal polynomial

$$\mu_A(x) = \prod_{i=1}^k (x - \lambda_i)^{d(\lambda_i)}.$$

This is just Theorem 2.7.4, (ii) together with Theorem 2.9.1, (ii).

3 Functions of matrices

3.1 Powers of matrices

The theory of Jordan canonical form we developed can be used to compute powers of matrices efficiently. Suppose we need to compute A^{2022} (and please appreciate the joke with the exponent constantly equalling the current year here and in the exercises...) where $A =$

$$\begin{pmatrix} -2 & 0 & 0 & 0 \\ 0 & -2 & 1 & 0 \\ 0 & 0 & -2 & 0 \\ 1 & 0 & -2 & -2 \end{pmatrix}$$

is the matrix from Example 10 in 2.10.

There are two practical ways of computing A^n by hand for a general matrix A and a very large n . The first one involves the JCF of A .

If $J = P^{-1}AP$ is the JCF of A then it is sufficient to compute J^n because of the telescoping product:

$$A^n = (PJP^{-1})^n = PJP^{-1}PJP^{-1}P \dots JP^{-1} = PJ^nP^{-1}.$$

How do we work out what J^n is? Firstly, we need to convince ourselves that

$$(B \oplus C)^n = B^n \oplus C^n$$

for square matrices B, C . We leave this as an exercise in understanding the multiplication of direct sums of matrices (it might help to look at some small examples!) and we have already required this when thinking about the minimal polynomial of direct sums of matrices. Clearly, it extends to the direct sum of any finite number of square matrices.

So we are left to consider what the power of an individual Jordan block is. Again a small example will help us:

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}^2 = \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix}, \dots, \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}^n = \begin{pmatrix} 1 & n \\ 0 & 1 \end{pmatrix}.$$

The eigenvalue being 1 hides things a little so let's do a slightly more complicated example.

$$\begin{pmatrix} 2 & 1 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \end{pmatrix}^2 = \begin{pmatrix} 4 & 4 & 1 \\ 0 & 4 & 4 \\ 0 & 0 & 4 \end{pmatrix}, \quad \begin{pmatrix} 2 & 1 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \end{pmatrix}^3 = \begin{pmatrix} 8 & 12 & 6 \\ 0 & 8 & 12 \\ 0 & 0 & 8 \end{pmatrix}.$$

At this point you should be willing to believe the following formula, which is left as an exercise (use induction!) to prove.

$$J_{\lambda,k}^n = \begin{pmatrix} \lambda^n & n\lambda^{n-1} & \dots & \binom{n}{k-2}\lambda^{n-k+2} & \binom{n}{k-1}\lambda^{n-k+1} \\ 0 & \lambda^n & \dots & \binom{n}{k-3}\lambda^{n-k+3} & \binom{n}{k-2}\lambda^{n-k+2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & \lambda^n & n\lambda^{n-1} \\ 0 & 0 & \dots & 0 & \lambda^n \end{pmatrix} \quad (6)$$

3 Functions of matrices

where $\binom{n}{t} = \frac{n!}{(n-t)!t!}$ is the choose-function (or binomial coefficient), interpreted as $\binom{n}{t} = 0$ whenever $t > n$.

Let us apply it to the matrix A above:

$$\begin{aligned} A^n = PJ^nP^{-1} &= \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & -2 & 0 \end{pmatrix} \begin{pmatrix} -2 & 1 & 0 & 0 \\ 0 & -2 & 0 & 0 \\ 0 & 0 & -2 & 1 \\ 0 & 0 & 0 & -2 \end{pmatrix}^n \begin{pmatrix} 0 & 2 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} = \\ & \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & -2 & 0 \end{pmatrix} \begin{pmatrix} (-2)^n & n(-2)^{n-1} & 0 & 0 \\ 0 & (-2)^n & 0 & 0 \\ 0 & 0 & (-2)^n & n(-2)^{n-1} \\ 0 & 0 & 0 & (-2)^n \end{pmatrix} \begin{pmatrix} 0 & 2 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix} = \\ & \begin{pmatrix} (-2)^n & 0 & 0 & 0 \\ 0 & (-2)^n & n(-2)^{n-1} & 0 \\ 0 & 0 & (-2)^n & 0 \\ n(-2)^{n-1} & 0 & n(-2)^n & (-2)^n \end{pmatrix}. \end{aligned}$$

The second method of computing A^n uses Lagrange's interpolation polynomial. It is less labour intensive and more suitable for pen-and-paper calculations.

Suppose $\psi(M) = 0$ for a polynomial $\psi(z)$, in practice we will choose $\psi(z)$ to be either the minimal or characteristic polynomial. Dividing with remainder gives $z^n = q(z)\psi(z) + h(z)$, and we conclude that

$$A^n = q(A)\psi(A) + h(A) = h(A).$$

Division with remainder may appear problematic² for large n but there is a shortcut. If we know the roots of $\psi(z)$, say $\alpha_1, \dots, \alpha_k$ with their multiplicities m_1, \dots, m_k , then $h(z)$ can be found by solving the system of simultaneous equations in coefficients of $h(z)$:

$$f^{(t)}(\alpha_j) = h^{(t)}(\alpha_j), \quad 1 \leq j \leq k, \quad 0 \leq t < m_j$$

where $f(z) = z^n$ and $f^{(t)}$ is the t -th derivative of f with respect to z . In other words, $h(z)$ is what is known as Lagrange's interpolation polynomial for the function z^n at the roots of $\psi(z)$. Note that we only ever need to take $h(z)$ to be a polynomial of degree $m_1 + \dots + m_k - 1$.

Let's use this to find A^n again for A as above. We know the minimal polynomial $\mu_A(z) = (z+2)^2$. Given $\mu_A(z)$ is degree 2 we can take the Lagrange interpolation of z^n at the roots of $(z+2)^2$ to be $h(z) = \alpha z + \beta$. To determine α and β we have to solve

$$\begin{cases} (-2)^n & = & h(-2) & = & -2\alpha + \beta \\ n(-2)^{n-1} & = & h'(-2) & = & \alpha \end{cases}$$

Solving them gives $\alpha = n(-2)^{n-1}$ and $\beta = (1-n)(-2)^n$. It follows that

$$A^n = n(-2)^{n-1}A + (1-n)(-2)^nI = \begin{pmatrix} (-2)^n & 0 & 0 & 0 \\ 0 & (-2)^n & n(-2)^{n-1} & 0 \\ 0 & 0 & (-2)^n & 0 \\ n(-2)^{n-1} & 0 & n(-2)^n & (-2)^n \end{pmatrix}.$$

²Try to divide z^{2022} by $z^2 + z + 1$ without reading any further.

3.2 Applications to difference equations

Let us consider an *initial value problem* for an *autonomous* system with discrete time:

$$\mathbf{x}(n+1) = A\mathbf{x}(n), \quad n \in \mathbb{N}, \quad \mathbf{x}(0) = w.$$

Here $\mathbf{x}(n) \in K^m$ is a sequence of vectors in a vector space over a field K . One thinks of $\mathbf{x}(n)$ as a state of the system at time n . The initial state is $\mathbf{x}(0) = w$. The $n \times n$ -matrix A with coefficients in K describes the evolution of the system. The adjective *autonomous* means that the evolution equation does not change with the time³.

It takes longer to formulate this problem than to solve it. The solution is straightforward:

$$\mathbf{x}(n) = A\mathbf{x}(n-1) = A^2\mathbf{x}(n-2) = \dots = A^n\mathbf{x}(0) = A^n w. \quad (7)$$

As a working example, let us consider the Fibonacci numbers:

$$F_0 = 0, \quad F_1 = 1 \quad \text{and} \quad F_n = F_{n-1} + F_{n-2} \quad (n \geq 2).$$

The recursion relations for them turn into

$$\begin{pmatrix} F_n \\ F_{n+1} \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} F_{n-1} \\ F_n \end{pmatrix}$$

so that (7) immediately yields a general solution

$$\begin{pmatrix} F_n \\ F_{n+1} \end{pmatrix} = A^n \begin{pmatrix} 0 \\ 1 \end{pmatrix} \quad \text{where} \quad A = \begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix}. \quad (8)$$

We compute the characteristic polynomial of A to be $c_A(z) = z^2 - z - 1$. Its discriminant is 5. The roots of $c_A(z)$ are the golden ratio $\lambda = (1 + \sqrt{5})/2$ and $1 - \lambda = (1 - \sqrt{5})/2$. It is useful to observe that

$$2\lambda - 1 = \sqrt{5} \quad \text{and} \quad \lambda(1 - \lambda) = -1.$$

Let us introduce the number $\mu_n = \lambda^n - (1 - \lambda)^n$. Suppose the Lagrange interpolation of z^n at the roots of $z^2 - z - 1$ is $h(z) = \alpha z + \beta$. The condition on the coefficients is given by

$$\begin{cases} \lambda^n & = & h(\lambda) & = & \alpha\lambda + \beta \\ (1 - \lambda)^n & = & h(1 - \lambda) & = & \alpha(1 - \lambda) + \beta \end{cases}$$

Solving them gives

$$\alpha = \mu_n / \sqrt{5} \quad \text{and} \quad \beta = \mu_{n-1} / \sqrt{5}.$$

It follows that

$$A^n = \alpha A + \beta = \mu_n / \sqrt{5} A + \mu_{n-1} / \sqrt{5} I_2 = \begin{pmatrix} \mu_{n-1} / \sqrt{5} & \mu_n / \sqrt{5} \\ \mu_n / \sqrt{5} & (\mu_n + \mu_{n-1}) / \sqrt{5} \end{pmatrix}.$$

³A nonautonomous system would be described by $\mathbf{x}(n+1) = A(n)\mathbf{x}(n)$ here.

3 Functions of matrices

Equation (8) immediately implies that

$$F_n = \mu_n / \sqrt{5} \text{ and } A^n = \begin{pmatrix} F_{n-1} & F_n \\ F_n & F_{n+1} \end{pmatrix} .$$

If we try and do this for more complicated difference equations, we could meet matrices which aren't diagonalisable. Here's an example (taken from the book by Kaye and Wilson, §14.11), done using Jordan canonical form.

Example. Let x_n, y_n, z_n be sequences of complex numbers satisfying

$$\begin{cases} x_{n+1} = 3x_n + z_n, \\ y_{n+1} = -x_n + y_n - z_n, \\ z_{n+1} = y_n + 2z_n. \end{cases}$$

with $x_0 = y_0 = z_0 = 1$.

We can write this as

$$\mathbf{v}_{n+1} = \begin{pmatrix} 3 & 0 & 1 \\ -1 & 1 & -1 \\ 0 & 1 & 2 \end{pmatrix} \mathbf{v}_n.$$

So we have

$$\mathbf{v}_n = A^n \mathbf{v}_0 = A^n \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

where A is the 3×3 matrix above.

We find that the JCF of A is $J = P^{-1}DP$ where

$$J = J_{2,3} = \begin{pmatrix} 2 & 1 & 0 \\ 0 & 2 & 1 \\ 0 & 0 & 2 \end{pmatrix}, \quad P = \begin{pmatrix} 1 & 1 & 1 \\ 0 & -1 & 0 \\ -1 & 0 & 0 \end{pmatrix}.$$

The formula for the entries of J^k for J a Jordan block tells us that

$$\begin{aligned} J^n &= \begin{pmatrix} 2^n & n2^{n-1} & \binom{n}{2}2^{n-2} \\ 0 & 2^n & n2^{n-1} \\ 0 & 0 & 2^n \end{pmatrix} \\ &= 2^n \begin{pmatrix} 1 & \frac{1}{2}n & \frac{1}{4}\binom{n}{2} \\ 0 & 1 & \frac{1}{2}n \\ 0 & 0 & 1 \end{pmatrix} \end{aligned}$$

3 Functions of matrices

We therefore have

$$\begin{aligned}
 A^n &= PJ^nP^{-1} \\
 &= 2^n \begin{pmatrix} 1 & 1 & 1 \\ 0 & -1 & 0 \\ -1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 & \frac{1}{2}n & \frac{1}{4}\binom{n}{2} \\ 0 & 1 & \frac{1}{2}n \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 0 & 0 & -1 \\ 0 & -1 & 0 \\ 1 & 1 & 1 \end{pmatrix} \\
 &= 2^n \begin{pmatrix} 1 & 1 + \frac{1}{2}n & 1 + \frac{1}{2}n + \frac{1}{4}\binom{n}{2} \\ 0 & -1 & -\frac{1}{2}n \\ -1 & -\frac{1}{2}n & -\frac{1}{4}\binom{n}{2} \end{pmatrix} \begin{pmatrix} 0 & 0 & -1 \\ 0 & -1 & 0 \\ 1 & 1 & 1 \end{pmatrix} \\
 &= 2^n \begin{pmatrix} 1 + \frac{1}{2}n + \frac{1}{4}\binom{n}{2} & \frac{1}{4}\binom{n}{2} & \frac{1}{2}n + \frac{1}{4}\binom{n}{2} \\ -\frac{1}{2}n & 1 - \frac{1}{2}n & -\frac{1}{2}n \\ -\frac{1}{4}\binom{n}{2} & \frac{1}{2}n - \frac{1}{4}\binom{n}{2} & 1 - \frac{1}{4}\binom{n}{2} \end{pmatrix}
 \end{aligned}$$

Finally, we obtain

$$A^n \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = 2^n \begin{pmatrix} 1 + n + \frac{3}{4}\binom{n}{2} \\ 1 - \frac{3}{2}n \\ 1 + \frac{1}{2}n - \frac{3}{4}\binom{n}{2} \end{pmatrix}$$

or equivalently, using the fact that $\binom{n}{2} = \frac{n(n-1)}{2}$,

$$\begin{cases} x_n &= 2^n \left(\frac{3}{4}n^2 + \frac{5}{8}n + 1 \right), \\ y_n &= 2^n \left(1 - \frac{3}{2}n \right), \\ z_n &= 2^n \left(-\frac{3}{4}n^2 + \frac{7}{8}n + 1 \right). \end{cases}$$

3.3 Motivation: Systems of Differential Equations

Suppose we want to expand our repertoire and solve a system of first-order simultaneous differential equations, say

$$\begin{aligned}
 \frac{da}{dt} &= 3a - 4b + 8c, \\
 \frac{db}{dt} &= a - c, \\
 \frac{dc}{dt} &= a + b + c.
 \end{aligned}$$

These are common in the Differential Equations course last year. Let's write the system in a different form. We consider $\mathbf{v}(t) = \begin{pmatrix} a(t) \\ b(t) \\ c(t) \end{pmatrix}$, a vector-valued function of time, and write the above system as

$$\frac{d\mathbf{v}}{dt} = A\mathbf{v}$$

3 Functions of matrices

where A is the matrix

$$\begin{pmatrix} 3 & -4 & 8 \\ 1 & 0 & -1 \\ 1 & 1 & 1 \end{pmatrix}.$$

“Aha!” we say. “We know the solution is $\mathbf{v}(t) = e^{tA}\mathbf{v}(0)$!” But then we pause, and say “Hang on, what does e^{tA} actually mean?” In the next section, we’ll use what we now know about special forms of matrices to define e^{tA} , and other functions of a matrix, in a sensible way that will make this actually work; and having got our definition, we’ll work out how to calculate with it.

3.4 Definition of a function of a matrix

Suppose we have a “nice” one variable complex-valued function $f(z)$. What is $f(A)$? In general, there is no natural answer. We had one for $f(z) = z^n$ in Section 3.1 and we choose to generalise this to define $f(A)$ using the Jordan canonical form of A as follows. Let $J = P^{-1}AP$ with $J = J_{\lambda_1, k_1} \oplus \cdots \oplus J_{\lambda_t, k_t}$ being the JCF of A . We define

$$f(A) = Pf(J)P^{-1}, \text{ where } f(J) = f(J_{\lambda_1, k_1}) \oplus \cdots \oplus f(J_{\lambda_t, k_t}),$$

and

$$f(J_{\lambda, k}) = \begin{pmatrix} f(\lambda) & f'(\lambda) & \cdots & f^{[k-2]}(\lambda) & f^{[k-1]}(\lambda) \\ 0 & f(\lambda) & \cdots & f^{[k-3]}(\lambda) & f^{[k-2]}(\lambda) \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & f(\lambda) & f'(\lambda) \\ 0 & 0 & \cdots & 0 & f(\lambda) \end{pmatrix}. \quad (9)$$

The notation $f^{[k]}(z)$ is known as the divided power derivative and defined as

$$f^{[k]}(z) := \frac{1}{k!} f^{(k)}(z).$$

So $f^{[1]} = f'$, $f^{[2]} = \frac{1}{2}f''$, $f^{[3]} = \frac{1}{6}f'''$, etc. As you might imagine, deciding exactly what a “nice” function is, and whether this is definition is sensible for functions defined by power series etc. is more analysis than it is algebra. Thus, in this course we will ignore such issues. We are mainly interested in the exponential of a matrix. Taylor’s series at zero of the exponential function is $\sum_{k=0}^{\infty} \frac{z^k}{k!}$ and so we might think that the following equation should be true.

$$e^A = I_n + A + \frac{A^2}{2} + \frac{A^3}{6} + \cdots = \sum_{k=0}^{\infty} \frac{A^k}{k!}. \quad (10)$$

It is indeed true, i.e. this coincides with our definition of $e^A = f(A)$ where f is the standard exponential function. Note, however, that not everything we know about the exponential function of complex numbers is true when we apply it to matrices. For example, it is *not* true that $e^{B+C} = e^B e^C$ for general matrices B and C ; you may wish to find an explicit example.

Let’s start by calculating e^A for a matrix A .

3 Functions of matrices

Example 11. Consider $A = \begin{pmatrix} 1 & 4 \\ 1 & 1 \end{pmatrix}$. This was Example 3 from Section 2.8 above, and we saw that $P^{-1}AP = J$ where

$$P = \begin{pmatrix} 2 & -2 \\ 1 & 1 \end{pmatrix}, \quad J = \begin{pmatrix} 3 & 0 \\ 0 & -1 \end{pmatrix}.$$

Hence

$$\begin{aligned} e^A &= Pe^JP^{-1} \\ &= \begin{pmatrix} 2 & -2 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} e^3 & 0 \\ 0 & e^{-1} \end{pmatrix} \begin{pmatrix} 2 & -2 \\ 1 & 1 \end{pmatrix}^{-1} \\ &= \frac{1}{4} \begin{pmatrix} 2 & -2 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} e^3 & 0 \\ 0 & e^{-1} \end{pmatrix} \begin{pmatrix} 1 & 2 \\ -1 & 2 \end{pmatrix} \\ &= \begin{pmatrix} \frac{1}{2}e^3 + \frac{1}{2}e^{-1} & e^3 - e^{-1} \\ \frac{1}{4}e^3 - \frac{1}{4}e^{-1} & \frac{1}{2}e^3 + \frac{1}{2}e^{-1} \end{pmatrix}. \end{aligned}$$

Let's see another way to calculate e^A . We can again use Lagrange's interpolation method, which is often easier in practice.

Example 12. We compute e^A for the matrix A from Example 10, Section 2.10, using Lagrange interpolation. Suppose that $h(z) = \alpha z + \beta$ is the interpolation of e^z at the roots of $\mu_A(z) = (z + 2)^2$. The condition on the coefficients is given by

$$\begin{cases} e^{-2} = h(-2) = -2\alpha + \beta \\ e^{-2} = h'(-2) = \alpha \end{cases}$$

Solving them gives $\alpha = e^{-2}$ and $\beta = 3e^{-2}$. It follows that

$$e^A = h(A) = e^{-2}A + 3e^{-2}I = \begin{pmatrix} e^{-2} & 0 & 0 & 0 \\ 0 & e^{-2} & e^{-2} & 0 \\ 0 & 0 & e^{-2} & 0 \\ e^{-2} & 0 & -2e^{-2} & e^{-2} \end{pmatrix}.$$

Our motivation for defining the exponential of a matrix was to find e^{tA} so let's do that in the next example. It is important to note that t here should be seen as a constant when we differentiate $f(z) = e^{zt}$. So $f^{[1]}(z) = te^{zt}$, $f^{[2]}(z) = \frac{1}{2}t^2e^{zt}$, etc.

Example 13. Let

$$A = \begin{pmatrix} 1 & 0 & -3 \\ 1 & -1 & -6 \\ -1 & 2 & 5 \end{pmatrix}.$$

Using the methods of the last chapter we can check that its JCF is $J = \begin{pmatrix} 2 & 1 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 1 \end{pmatrix}$ and the basis

change matrix P such that $J = P^{-1}AP$ is given by $P = \begin{pmatrix} 3 & 0 & 2 \\ 3 & 1 & 1 \\ -1 & -1 & 0 \end{pmatrix}$.

3 Functions of matrices

Applying the argument above, we see that $e^{tA} = Pe^{tJ}P^{-1}$ where

$$e^{tJ} = \begin{pmatrix} e^{2t} & te^{2t} & 0 \\ 0 & e^{2t} & 0 \\ 0 & 0 & e^t \end{pmatrix}.$$

We can now calculate e^{tA} explicitly by doing the matrix multiplication to get the entries of $Pe^{tJ}P^{-1}$, as we did in the 2×2 example above.

It looks messy. Do we really want to write it down here?

Well, let us not do it. In a pen-and-paper calculation, except a few cases (for example, diagonal matrices) it is simpler to use Lagrange's interpolation.

Example 14. Let us consider a harmonic oscillator described by the equation $y''(t) + y(t) = 0$. The general solution $y(t) = \alpha \sin(t) + \beta \cos(t)$ is well known. Let us obtain it using matrix exponents. Setting

$$x(t) = \begin{pmatrix} y(t) \\ y'(t) \end{pmatrix}, \quad A = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$$

the harmonic oscillator becomes the initial value problem with a solution $x(t) = e^{tA}x(0)$. The eigenvalues of A are i and $-i$. Interpolating e^{tz} at these values of z gives the following condition on $h(z) = \alpha z + \beta$

$$\begin{cases} e^{ti} & = & h(i) & = & \alpha i + \beta \\ e^{-ti} & = & h(-i) & = & -\alpha i + \beta \end{cases}$$

Solving them gives $\alpha = (e^{ti} - e^{-ti})/2i = \sin(t)$ and $\beta = (e^{ti} + e^{-ti})/2 = \cos(t)$. It follows that

$$e^{tA} = \sin(t)A + \cos(t)I_2 = \begin{pmatrix} \cos(t) & \sin(t) \\ -\sin(t) & \cos(t) \end{pmatrix}$$

and so

$$x(t) = \begin{pmatrix} \cos(t)y(0) + \sin(t)y'(0) \\ -\sin(t)y(0) + \cos(t)y'(0) \end{pmatrix}.$$

The final solution is thus $y(t) = \cos(t)y(0) + \sin(t)y'(0)$.

Example 15. Let us consider a system of differential equations

$$\begin{cases} y_1' & = & y_1 - 3y_3 \\ y_2' & = & y_1 - y_2 - 6y_3 \\ y_3' & = & -y_1 + 2y_2 + 5y_3 \end{cases}, \quad \text{with the initial condition } \begin{cases} y_1(0) & = & 1 \\ y_2(0) & = & 1 \\ y_3(0) & = & 0 \end{cases}$$

Using matrices

$$x(t) = \begin{pmatrix} y_1(t) \\ y_2(t) \\ y_3(t) \end{pmatrix}, \quad w = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}, \quad A = \begin{pmatrix} 1 & 0 & -3 \\ 1 & -1 & -6 \\ -1 & 2 & 5 \end{pmatrix},$$

3 Functions of matrices

it becomes an initial value problem. The characteristic polynomial is $c_A(z) = -z^3 + 5z^2 - 8z + 4 = (1-z)(2-z)^2$. We need to interpolate e^{tz} at 1 and 2 by $h(z) = \alpha z^2 + \beta z + \gamma$. At the multiple root 2 we need to interpolate up to order 2 that involves tracking the derivative $(e^{tz})' = te^{tz}$:

$$\begin{cases} e^t &= h(1) &= \alpha + \beta + \gamma \\ e^{2t} &= h(2) &= 4\alpha + 2\beta + \gamma \\ te^{2t} &= h'(2) &= 4\alpha + \beta \end{cases}$$

Solving, $\alpha = (t-1)e^{2t} + e^t$, $\beta = (4-3t)e^{2t} - 4e^t$, $\gamma = (2t-3)e^{2t} + 4e^t$. It follows that

$$e^{tA} = e^{2t} \begin{pmatrix} 3t-3 & -6t+6 & -9t+6 \\ 3t-2 & -6t+4 & -9t+3 \\ -t & 2t & 3t+1 \end{pmatrix} + e^t \begin{pmatrix} 4 & -6 & -6 \\ 2 & -3 & -3 \\ 0 & 0 & 0 \end{pmatrix}$$

and

$$x(t) = \begin{pmatrix} y_1(t) \\ y_2(t) \\ y_3(t) \end{pmatrix} = e^{tA} \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} (3-3t)e^{2t} - 2e^t \\ (2-3t)e^{2t} - e^t \\ te^{2t} \end{pmatrix}.$$

4 Bilinear Maps and Quadratic Forms

We'll now introduce another, rather different kind of object you can define for vector spaces: a *bilinear map*. These are a bit different from linear maps: rather than being machines that take a vector and spit out another vector, they take two vectors as input and spit out a number.

4.1 Bilinear maps: definitions

Let V and W be vector spaces over a field K .

Definition 4.1.1. A *bilinear map* on V and W is a map $\tau : V \times W \rightarrow K$ such that

- (i) $\tau(\alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2, \mathbf{w}) = \alpha_1 \tau(\mathbf{v}_1, \mathbf{w}) + \alpha_2 \tau(\mathbf{v}_2, \mathbf{w})$; and
- (ii) $\tau(\mathbf{v}, \alpha_1 \mathbf{w}_1 + \alpha_2 \mathbf{w}_2) = \alpha_1 \tau(\mathbf{v}, \mathbf{w}_1) + \alpha_2 \tau(\mathbf{v}, \mathbf{w}_2)$

for all $\mathbf{v}, \mathbf{v}_1, \mathbf{v}_2 \in V$, $\mathbf{w}, \mathbf{w}_1, \mathbf{w}_2 \in W$, and $\alpha_1, \alpha_2 \in K$.

So $\tau(\mathbf{v}, \mathbf{w})$ is linear in \mathbf{v} for each \mathbf{w} , and linear in \mathbf{w} for each \mathbf{v} – linear in two different ways, hence the term “bilinear”.

Clearly if we fix bases of V and W , a bilinear map will be determined by what it does to the basis vectors. Choose a basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ of V and a basis $\mathbf{f}_1, \dots, \mathbf{f}_m$ of W .

Let $\tau : V \times W \rightarrow K$ be a bilinear map, and let $\alpha_{ij} = \tau(\mathbf{e}_i, \mathbf{f}_j)$, for $1 \leq i \leq n, 1 \leq j \leq m$. Then the $n \times m$ matrix $A = (\alpha_{ij})$ is defined to be the matrix of τ with respect to the bases $\mathbf{e}_1, \dots, \mathbf{e}_n$ and $\mathbf{f}_1, \dots, \mathbf{f}_m$ of V and W .

For $\mathbf{v} \in V, \mathbf{w} \in W$, let $\mathbf{v} = x_1 \mathbf{e}_1 + \dots + x_n \mathbf{e}_n$ and $\mathbf{w} = y_1 \mathbf{f}_1 + \dots + y_m \mathbf{f}_m$, so the coordinates of \mathbf{v} and \mathbf{w} with respect to our bases are

$$\underline{\mathbf{v}} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in K^{n,1}, \quad \text{and} \quad \underline{\mathbf{w}} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} \in K^{m,1}.$$

Then, by using the equations (i) and (ii) above, we get

$$\tau(\mathbf{v}, \mathbf{w}) = \sum_{i=1}^n \sum_{j=1}^m x_i \tau(\mathbf{e}_i, \mathbf{f}_j) y_j = \sum_{i=1}^n \sum_{j=1}^m x_i \alpha_{ij} y_j = \underline{\mathbf{v}}^T A \underline{\mathbf{w}}. \quad (2.1)$$

So once we've fixed bases of V and W , every bilinear map on V and W corresponds to an $n \times m$ matrix, and conversely every matrix determines a bilinear map.

For example, let $V = W = \mathbb{R}^2$ and use the natural basis of V . Suppose that $A = \begin{pmatrix} 1 & -1 \\ 2 & 0 \end{pmatrix}$.

Then

$$\tau((x_1, x_2), (y_1, y_2)) = (x_1 \ x_2) \begin{pmatrix} 1 & -1 \\ 2 & 0 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = x_1 y_1 - x_1 y_2 + 2x_2 y_1.$$

4.2 Bilinear maps: change of basis

We retain the notation of the previous section, so τ is a bilinear map on V and W , and A is the matrix of τ with respect to some bases $\mathbf{e}_1, \dots, \mathbf{e}_n$ of V and $\mathbf{f}_1, \dots, \mathbf{f}_m$ of W .

As in §1.5 of the course, suppose that we choose new bases $\mathbf{e}'_1, \dots, \mathbf{e}'_n$ of V and $\mathbf{f}'_1, \dots, \mathbf{f}'_m$ of W , and let P and Q be the associated basis change matrices. Let B be the matrix of τ with respect to these new bases.

Let \mathbf{v} be any vector in V . Then we know (from Proposition 1.5.1) that if $\underline{\mathbf{v}} \in K^{n,1}$ is the column vector of coordinates of \mathbf{v} with respect to the old basis $\mathbf{e}_1, \dots, \mathbf{e}_n$, and $\underline{\mathbf{v}}'$ the coordinates of \mathbf{v} in the new basis $\mathbf{e}'_1, \dots, \mathbf{e}'_n$, then we have $P\underline{\mathbf{v}}' = \underline{\mathbf{v}}$. Similarly, for any $\mathbf{w} \in W$, the coordinates $\underline{\mathbf{w}}$ and $\underline{\mathbf{w}}'$ of \mathbf{w} with respect to the old and new bases of W are related by $Q\underline{\mathbf{w}}' = \underline{\mathbf{w}}$.

We know that we have

$$\underline{\mathbf{v}}^T A \underline{\mathbf{w}} = \tau(\mathbf{v}, \mathbf{w}) = (\underline{\mathbf{v}}')^T B \underline{\mathbf{w}}'.$$

Substituting in the formulae from Proposition 1.5.1, we have

$$\begin{aligned} (\underline{\mathbf{v}}')^T B \underline{\mathbf{w}}' &= (P\underline{\mathbf{v}}')^T A (Q\underline{\mathbf{w}}') \\ &= (\underline{\mathbf{v}}')^T P^T A Q \underline{\mathbf{w}}'. \end{aligned}$$

Since this relation must hold for all $\underline{\mathbf{v}}' \in K^{n,1}$ and $\underline{\mathbf{w}}' \in K^{m,1}$, the two matrices in the middle must be equal (exercise!): that is, we have $B = P^T A Q$. So we've proven:

Theorem 4.2.1. *Let A be the matrix of the bilinear map $\tau : V \times W \rightarrow K$ with respect to the bases $\mathbf{e}_1, \dots, \mathbf{e}_n$ and $\mathbf{f}_1, \dots, \mathbf{f}_m$ of V and W , and let B be its matrix with respect to the bases $\mathbf{e}'_1, \dots, \mathbf{e}'_n$ and $\mathbf{f}'_1, \dots, \mathbf{f}'_m$ of V and W . Let P and Q be the basis change matrices, as defined above. Then $B = P^T A Q$.*

Compare this result with Theorem 1.5.2.

We shall be concerned from now on only with the case where $V = W$. A bilinear map $\tau : V \times V \rightarrow K$ is called a *bilinear form* on V . Theorem 4.2.1 then becomes:

Theorem 4.2.2. *Let A be the matrix of the bilinear form τ on V with respect to the basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ of V , and let B be its matrix with respect to the basis $\mathbf{e}'_1, \dots, \mathbf{e}'_n$ of V . Let P the basis change matrix with original basis $\{\mathbf{e}_i\}$ and new basis $\{\mathbf{e}'_i\}$. Then $B = P^T A P$.*

Let's give a name to this relation between matrices:

Definition 4.2.3. Two matrices A and B are called *congruent* if there exists an invertible matrix P with $B = P^T A P$.

So congruent matrices represent the same bilinear form in different bases. Notice that congruence is very different from similarity; if τ is a bilinear form on V and T is a linear operator on V , it might be the case that τ and T have the same matrix A in some specific basis of V , but that doesn't mean that they have the same matrix in any other basis – they inhabit different worlds.

4 Bilinear Maps and Quadratic Forms

So, in the example at the end of Subsection 4.1, if we choose the new basis $\mathbf{e}'_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$, $\mathbf{e}'_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ then $P = \begin{pmatrix} 1 & 1 \\ -1 & 0 \end{pmatrix}$, $P^T A P = \begin{pmatrix} 0 & -1 \\ 2 & 1 \end{pmatrix}$, and

$$\tau((y'_1 \mathbf{e}'_1 + y'_2 \mathbf{e}'_2, x'_1 \mathbf{e}'_1 + x'_2 \mathbf{e}'_2)) = -y'_1 x'_2 + 2y'_2 x'_1 + y'_2 x'_2.$$

Since P is an invertible matrix, P^T is also invertible (its inverse is $(P^{-1})^T$), and so the matrices $P^T A P$ and A are “equivalent matrices” in the sense of MA106, and hence have the same rank.

The rank of the bilinear form τ is defined to be the rank of its matrix A . So we have just shown that the rank of τ is a well-defined property of τ , not depending on the choice of basis we’ve used.

In fact we can say a little more. It’s clear that a vector $\underline{\mathbf{v}} \in K^{n,1}$ is zero if and only if $\underline{\mathbf{v}}^T \underline{\mathbf{w}} = 0$ for all vectors $\underline{\mathbf{w}} \in K^{n,1}$. Since

$$\tau(\mathbf{v}, \mathbf{w}) = \underline{\mathbf{v}}^T A \underline{\mathbf{w}},$$

the kernel of A is equal to the space

$$\{\mathbf{v} \in V : \tau(\mathbf{w}, \mathbf{v}) = 0 \forall \mathbf{w} \in V\}$$

(the *right radical* of τ) and the kernel of A^T is equal to the space

$$\{\mathbf{v} \in V : \tau(\mathbf{v}, \mathbf{w}) = 0 \forall \mathbf{w} \in V\}$$

(the *left radical*). Since A^T and A have the same rank, the left and right radicals both have dimension $n - r$, where r is the rank of τ . In particular, the rank of τ is n if and only if the left and right radicals are zero. If this occurs, we’ll say τ is *nondegenerate*; so τ is nondegenerate if and only if its matrix (in any basis) is nonsingular.

You could be forgiven for expecting that we were about to launch into a long study of how to choose, given a bilinear form τ on V , the “best” basis for V which makes the matrix of τ as nice as possible. We are *not* going to do this, because although it’s a very natural question to ask, it’s *extremely* hard! Instead, we’ll restrict ourselves to a special kind of bilinear form where life is much easier, which covers most of the bilinear forms that come up in “real life”.

Definition 4.2.4. We say bilinear form τ on V is *symmetric* if $\tau(\mathbf{w}, \mathbf{v}) = \tau(\mathbf{v}, \mathbf{w})$ for all $\mathbf{v}, \mathbf{w} \in V$.

We say τ is *antisymmetric* (or sometimes *alternating*) if $\tau(\mathbf{v}, \mathbf{v}) = 0$ for all $\mathbf{v} \in V$.

The antisymmetry condition implies for all $\mathbf{v}, \mathbf{w} \in V$

$$\tau(\mathbf{v} + \mathbf{w}, \mathbf{v} + \mathbf{w}) = \tau(\mathbf{v}, \mathbf{w}) + \tau(\mathbf{w}, \mathbf{v}) = 0$$

hence for all $\mathbf{v}, \mathbf{w} \in V$

$$\tau(\mathbf{v}, \mathbf{w}) = -\tau(\mathbf{w}, \mathbf{v}).$$

If $2 \neq 0$ in K , the condition $\tau(\mathbf{v}, \mathbf{w}) = -\tau(\mathbf{w}, \mathbf{v})$ implies antisymmetry (take $\mathbf{v} = \mathbf{w}$, but you need to be able to divide by 2).

An $n \times n$ matrix A is called symmetric if $A^T = A$, and anti-symmetric if $A^T = -A$ and it has zeros along the diagonal. We then clearly have:

4 Bilinear Maps and Quadratic Forms

Proposition 4.2.5. *The bilinear form τ is symmetric or anti-symmetric if and only if its matrix (with respect to any basis) is symmetric or anti-symmetric.*

The best known example of a symmetric form is when $V = \mathbb{R}^n$, and τ is defined by

$$\tau \left(\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}, \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \right) = x_1y_1 + x_2y_2 + \cdots + x_ny_n.$$

This form has matrix equal to the identity matrix I_n with respect to the standard basis of \mathbb{R}^n . Geometrically, it is equal to the normal scalar product $\tau(\mathbf{v}, \mathbf{w}) = |\mathbf{v}||\mathbf{w}| \cos \theta$, where θ is the angle between the vectors \mathbf{v} and \mathbf{w} .

On the other hand, the form on \mathbb{R}^2 defined by $\tau \left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \right) = x_1y_2 - x_2y_1$ is anti-symmetric.

This has matrix $\begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$.

Proposition 4.2.6. *Suppose that $2 \neq 0$ in K . Then any bilinear form τ can be written uniquely as $\tau_1 + \tau_2$ where τ_1 is symmetric and τ_2 is antisymmetric.*

Proof. We just put $\tau_1(\mathbf{v}, \mathbf{w}) = \frac{1}{2}(\tau(\mathbf{v}, \mathbf{w}) + \tau(\mathbf{w}, \mathbf{v}))$ and $\tau_2(\mathbf{v}, \mathbf{w}) = \frac{1}{2}(\tau(\mathbf{v}, \mathbf{w}) - \tau(\mathbf{w}, \mathbf{v}))$. It's clear that τ_1 is symmetric and τ_2 is antisymmetric.

Moreover, given any other such expression $\tau = \tau'_1 + \tau'_2$, we have

$$\begin{aligned} \tau_1(\mathbf{v}, \mathbf{w}) &= \frac{\tau'_1(\mathbf{v}, \mathbf{w}) + \tau'_1(\mathbf{w}, \mathbf{v}) + \tau'_2(\mathbf{v}, \mathbf{w}) + \tau'_2(\mathbf{w}, \mathbf{v})}{2} \\ &= \frac{\tau'_1(\mathbf{v}, \mathbf{w}) + \tau'_1(\mathbf{v}, \mathbf{w}) + \tau'_2(\mathbf{v}, \mathbf{w}) - \tau'_2(\mathbf{v}, \mathbf{w})}{2} \end{aligned}$$

from the symmetry and antisymmetry of τ'_1 and τ'_2 . The last two terms cancel each other and we just have

$$= \frac{2\tau'_1(\mathbf{v}, \mathbf{w})}{2} = \tau'_1(\mathbf{v}, \mathbf{w}).$$

So $\tau_1 = \tau'_1$, and so $\tau_2 = \tau - \tau_1 = \tau - \tau'_1 = \tau'_2$, so the decomposition is unique. □

(Notice that $\frac{1}{2}$ has to exist in K for all this to make sense!)

4.3 Quadratic forms

Definition 4.3.1. Let V be a vector space over the field K . Then a *quadratic form* on V is a function $q : V \rightarrow K$ that satisfies that

$$q(\lambda\mathbf{v}) = \lambda^2q(\mathbf{v}), \quad \forall \mathbf{v} \in V, \lambda \in K$$

4 Bilinear Maps and Quadratic Forms

and that

$$(*) \quad \tau_q(\mathbf{v}_1, \mathbf{v}_2) := q(\mathbf{v}_1 + \mathbf{v}_2) - q(\mathbf{v}_1) - q(\mathbf{v}_2)$$

is a symmetric bilinear form on V .

As we can see from the definition, symmetric bilinear forms and quadratic forms are closely related. Indeed, given a bilinear form τ we can define a quadratic form by

$$q_\tau(\mathbf{v}) := \tau(\mathbf{v}, \mathbf{v}).$$

Moreover, given a quadratic form, $(*)$ above gives us a symmetric bilinear form. These processes are *almost* inverse to each other: indeed, one can easily compute that starting with a quadratic form q and bilinear form τ

$$q_{\tau_q} = 2q, \quad \tau_{q_\tau} = 2\tau.$$

So as long as $2 \neq 0$ in our K , quadratic forms and bilinear forms correspond to each other in a one-to-one way if we make the associations

$$q \mapsto \frac{1}{2}\tau_q, \quad \tau \mapsto q_\tau.$$

If $2 = 0$ in K (e.g. in $\mathbb{F}_2 = \mathbb{Z}/2\mathbb{Z}$, but there are again lots of other examples of such fields) this correspondence breaks down: indeed, in that case there are quadratic forms that are *not* of the form $\tau(-, -)$ for any symmetric bilinear form τ on V ; e.g. let $V = \mathbb{F}_2^2$, the space of pairs (x_1, x_2) with $x_i \in \mathbb{F}_2$. We would certainly like to be able to call

$$q((x_1, x_2)) = x_1x_2$$

a quadratic form on V . On the other hand, a general symmetric bilinear form on V looks like

$$\tau((x_1, x_2), (y_1, y_2)) = ax_1y_1 + bx_1y_2 + bx_2y_1 + cx_2y_2$$

so that putting $(x_1, x_2) = (y_1, y_2)$ we only get quadratic forms that are sums of squares.

There is an important and highly developed theory of quadratic forms also when $2 = 0$ in K (exposed in for example the books by Merkurjev-Karpenko-Elman or Kneser on quadratic forms), but the normal forms for them are a bit different from the case when $2 \neq 0$ and though the theory is not actually harder it divides naturally according to whether $2 = 0$ or $2 \neq 0$ in K . So from now on till the rest of this Chapter we make the:

Assumption: In our field K , we have that $2 = 1 + 1$ is not equal to 0.

Let $\mathbf{e}_1, \dots, \mathbf{e}_n$ be a basis of V . Recall that the coordinates of \mathbf{v} with respect to this basis are defined to be the field elements x_i such that $\mathbf{v} = \sum_{i=1}^n x_i \mathbf{e}_i$.

Let $A = (a_{ij})$ be the matrix of a symmetric bilinear form τ with respect to this basis. We will also call A the matrix of $q = q_\tau$ with respect to this basis. Then A is symmetric because τ is, and by Equation (2.1) of Subsection 4.1, we have

4 Bilinear Maps and Quadratic Forms

$$q(\mathbf{v}) = \mathbf{v}^T A \mathbf{v} = \sum_{i=1}^n \sum_{j=1}^n x_i \alpha_{ij} x_j = \sum_{i=1}^n \alpha_{ii} x_i^2 + 2 \sum_{i=1}^n \sum_{j=1}^{i-1} \alpha_{ij} x_i x_j. \quad (3.1)$$

Conversely, if we are given a quadratic form as in the right hand side of Equation (3.1), then it is easy to write down its matrix A . For example, if $n = 3$ and $q(\mathbf{v}) = 3x^2 + y^2 - 2z^2 + 4xy - xz$,

then $A = \begin{pmatrix} 3 & 2 & -1/2 \\ 2 & 1 & 0 \\ -1/2 & 0 & -2 \end{pmatrix}$.

4.4 Nice bases for quadratic forms

We'll now show how to choose a basis for V which makes a given symmetric bilinear form (or, equivalently, quadratic form) "as nice as possible". This will turn out to be much easier than the corresponding problem for linear operators.

Theorem 4.4.1. *Let V be a vector space of dimension n equipped with a symmetric bilinear form τ (or, equivalently, a quadratic form q).*

Then there is a basis $\mathbf{b}_1, \dots, \mathbf{b}_n$ of V , and constants β_1, \dots, β_n , such that

$$\tau(\mathbf{b}_i, \mathbf{b}_j) = \begin{cases} \beta_i & \text{if } j = i \\ 0 & \text{if } j \neq i \end{cases}$$

Equivalently,

- *given any symmetric matrix A , there is an invertible matrix P such that $P^T A P$ is a diagonal matrix (i.e. A is congruent to a diagonal matrix);*
- *given any quadratic form q on a vector space V , there is a basis $\mathbf{b}_1, \dots, \mathbf{b}_n$ of V and constants β_1, \dots, β_n such that*

$$q(x_1 \mathbf{b}_1 + \dots + x_n \mathbf{b}_n) = \beta_1 x_1^2 + \dots + \beta_n x_n^2.$$

Proof. We shall prove this by induction on $n = \dim V$. If $n = 0$ then there is nothing to prove, so let's assume that $n \geq 1$.

If τ is zero, then again there is nothing to prove, so we may assume that $\tau \neq 0$. Then the associated quadratic form q is not zero either, so there is a vector $\mathbf{v} \in V$ such that $q(\mathbf{v}, \mathbf{v}) \neq 0$. Let $\mathbf{b}_1 = \mathbf{v}$ and let $\beta_1 = q(\mathbf{v})$.

Consider the linear map $V \rightarrow K$ given by $\mathbf{w} \mapsto \tau(\mathbf{w}, \mathbf{v})$. This is not the zero map, so its image has rank 1; so its kernel W has rank $n - 1$. Moreover, this $(n - 1)$ -dimensional subspace doesn't contain $\mathbf{b}_1 = \mathbf{v}$.

By the induction hypothesis, we can find a basis $\mathbf{b}_2, \dots, \mathbf{b}_n$ for the kernel such that $\tau(\mathbf{b}_i, \mathbf{b}_j) = 0$ for all $2 \leq i < j \leq n$; and all of these vectors lie in the space W , so we also have $\tau(\mathbf{b}_1, \mathbf{b}_j) = 0$ for all $2 \leq j \leq n$. Since $\mathbf{b}_1 \notin W$, it follows that $\mathbf{b}_1, \dots, \mathbf{b}_n$ is a basis of V , so we're done. \square

4 Bilinear Maps and Quadratic Forms

Finding the good basis: The above proof is quite short and slick, and gives us very little help if we explicitly want to find the diagonalizing basis. So let's unravel what's going on a bit more explicitly. We'll see in a moment that what's going on is very closely related to "completing the square" in school algebra.

So let's say we have a quadratic form q . As usual, let $B = (\beta_{ij})$ be the matrix of q with respect to some arbitrary basis $\mathbf{b}_1, \dots, \mathbf{b}_n$. We'll modify the basis \mathbf{b}_i step-by-step in order to eventually get it into the nice form the theorem predicts.

Step 1: Arrange that $q(\mathbf{b}_1) \neq 0$. Here there are various cases to consider.

- If $\beta_{11} \neq 0$, then we're done: this means that $q(\mathbf{b}_1) \neq 0$, so we don't need to do anything.
- If $\beta_{11} = 0$, but $\beta_{ii} \neq 0$ for some $i > 1$, then we just interchange \mathbf{b}_1 and \mathbf{b}_i in our basis.
- If $\beta_{ii} = 0$ for all i , but there is some i and j such that $\beta_{ij} \neq 0$, then we replace \mathbf{b}_i with $\mathbf{b}_i + \mathbf{b}_j$; since

$$q(\mathbf{b}_i + \mathbf{b}_j) = q(\mathbf{b}_i) + q(\mathbf{b}_j) + 2\tau(\mathbf{b}_i, \mathbf{b}_j) = 2\beta_{ij},$$

after making this change we have $q(\mathbf{b}_i) \neq 0$, so we're reduced to one of the two previous cases.

- If $\beta_{ij} = 0$ for all i and j , we can stop: the matrix of q is zero, so it's certainly diagonal.

Step 2: Modify $\mathbf{b}_2, \dots, \mathbf{b}_n$ to make them orthogonal to \mathbf{b}_1 . Suppose we've done Step 1, but we haven't stopped, so β_{11} is now non-zero. We want to arrange that $\tau(\mathbf{b}_1, \mathbf{b}_i)$ is 0 for all $i > 1$. To do this, we just replace \mathbf{b}_i with

$$\mathbf{b}_i - \frac{\beta_{1i}}{\beta_{11}} \mathbf{b}_1.$$

This works because

$$\tau(\mathbf{b}_1, \mathbf{b}_i - \frac{\beta_{1i}}{\beta_{11}} \mathbf{b}_1) = \tau(\mathbf{b}_1, \mathbf{b}_i) - \frac{\beta_{1i}}{\beta_{11}} \tau(\mathbf{b}_1, \mathbf{b}_1) = \beta_{1i} - \frac{\beta_{1i}}{\beta_{11}} \beta_{11} = 0.$$

This is where the relation to "completing the square" comes in. We've changed our basis by the matrix

$$P = \begin{pmatrix} 1 & -\frac{\beta_{12}}{\beta_{11}} & \cdots & -\frac{\beta_{1n}}{\beta_{11}} \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix}$$

so the coordinates of a vector $\mathbf{v} \in V$ change by the inverse of this, which is just

$$P^{-1} = \begin{pmatrix} 1 & \frac{\beta_{12}}{\beta_{11}} & \cdots & \frac{\beta_{1n}}{\beta_{11}} \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix}$$

4 Bilinear Maps and Quadratic Forms

This corresponds to writing

$$q(x_1 \mathbf{b}_1 + \cdots + x_n \mathbf{b}_n) = \beta_{11} x_1^2 + 2\beta_{12} x_1 x_2 + \cdots + 2\beta_{1n} x_1 x_n + C$$

where C doesn't involve x_1 at all, and writing this as

$$\beta_{11} \left(x_1 + \frac{\beta_{12}}{\beta_{11}} x_2 + \cdots + \frac{\beta_{1n}}{\beta_{11}} x_n \right)^2 + C'$$

where C' also doesn't involve x_1 . Then our change of basis changes the coordinates so the whole bracketed term becomes the first coordinate of \mathbf{v} ; we've eliminated "cross terms" involving x_1 and one of the other variables.

Step 3: Induct on n . Now we've managed to engineer a basis $\mathbf{b}_1, \dots, \mathbf{b}_n$ such that the matrix $B = \beta_{ij}$ of q looks like

$$\begin{pmatrix} \beta_{11} & 0 & \cdots & 0 \\ 0 & ? & \cdots & ? \\ \vdots & \vdots & \ddots & \vdots \\ 0 & ? & \cdots & ? \end{pmatrix}$$

So we can now repeat the process with V replaced by the $(n-1)$ -dimensional vector space W spanned by $\mathbf{b}_2, \dots, \mathbf{b}_n$. We can mess around as much as we like with the vectors $\mathbf{b}_2, \dots, \mathbf{b}_n$ without breaking the fact that they pair to zero with \mathbf{b}_1 , since this is true of any vector in W . So we go back to step 1 but with a smaller n , and keep going until we either have a 0-dimensional space or a zero form, in which case we can safely stop.

Example. Let $V = \mathbb{R}^3$ and $q \left(\begin{pmatrix} x \\ y \\ z \end{pmatrix} \right) = xy + 3yz - 5xz$, so the matrix of q with respect to the standard basis $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$ is

$$A = \begin{pmatrix} 0 & 1/2 & -5/2 \\ 1/2 & 0 & 3/2 \\ -5/2 & 3/2 & 0 \end{pmatrix}.$$

Since we have only 3 variables, it's much less work to call them x, y, z than x_1, x_2, x_3 . When we change the variables, we will write x_1, y_1, z_1 and so on. We still proceed as in the previous proof and you need to read the proof first! We will use $\stackrel{\heartsuit}{=}$ for the equalities that need no checking (they are for information purposes only).

First change of basis: All the diagonal entries of A are zero, so we're in Case 3 of Step 1 of the proof above. But a_{12} is $1/2$, which isn't zero; so we replace \mathbf{e}_1 with $\mathbf{e}_1 + \mathbf{e}_2$. That is, we work in the basis

$$\mathbf{b}_1 := \mathbf{e}_1 + \mathbf{e}_2, \quad \mathbf{b}_2 := \mathbf{e}_2, \quad \mathbf{b}_3 := \mathbf{e}_3.$$

Thus the basis change matrix from $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$ to $\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3$ is

$$P = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{so that} \quad \begin{pmatrix} x \\ y \\ z \end{pmatrix} \stackrel{\heartsuit}{=} P \begin{pmatrix} x_1 \\ y_1 \\ z_1 \end{pmatrix}$$

4 Bilinear Maps and Quadratic Forms

where $\begin{pmatrix} x_1 \\ y_1 \\ z_1 \end{pmatrix}$ is the coordinate expression in the new basis (remember, P takes new coordinates to old). And we have

$$\begin{aligned} q(x_1\mathbf{b}_1 + y_1\mathbf{b}_2 + z_1\mathbf{b}_3) &= q\begin{pmatrix} x_1 \\ x_1 + y_1 \\ z_1 \end{pmatrix} = \\ &= x_1(x_1 + y_1) + 3(x_1 + y_1)z_1 - 5x_1z_1 = x_1^2 + x_1y_1 - 2x_1z_1 + 3y_1z_1, \end{aligned}$$

so the matrix of q in the basis $\mathbf{b}_1, \mathbf{b}_2, \mathbf{b}_3$ is

$$B = \begin{pmatrix} 1 & 1/2 & -1 \\ 1/2 & 0 & 3/2 \\ -1 & 3/2 & 0 \end{pmatrix} \stackrel{\heartsuit}{=} P^T A P.$$

Second change of basis: Now we can use Step 2 of the proof to clear the entries in the first row and column by modifying \mathbf{b}_2 and \mathbf{b}_3 , this is the “completing the square” step. As specified in Step 2 of the proof, we introduce a new basis \mathbf{b}' as follows

$$\mathbf{b}'_1 := \mathbf{b}_1 = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}, \quad \mathbf{b}'_2 := \mathbf{b}_2 - \frac{1}{2}\mathbf{b}_1 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} - \frac{1}{2}\begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} -1/2 \\ 1/2 \\ 0 \end{pmatrix},$$

$$\mathbf{b}'_3 := \mathbf{b}_3 - (-1)\mathbf{b}_1 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}.$$

So the basis change matrix from $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$ to $\mathbf{b}'_1, \mathbf{b}'_2, \mathbf{b}'_3$ is

$$P' = \begin{pmatrix} 1 & -1/2 & 1 \\ 1 & 1/2 & 1 \\ 0 & 0 & 1 \end{pmatrix} \stackrel{\heartsuit}{=} P Q \quad \text{where} \quad Q = \begin{pmatrix} 1 & -1/2 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

This corresponds to writing

$$\begin{aligned} [x_1^2 + x_1y_1 - 2x_1z_1] + 3y_1z_1 &= \left[(x_1 + \frac{1}{2}y_1 - z_1)^2 - \frac{1}{4}y_1^2 - z_1^2 + y_1z_1 \right] + 3y_1z_1 \\ &= (x_1 + \frac{1}{2}y_1 - z_1)^2 - \frac{1}{4}y_1^2 + 4y_1z_1 - z_1^2. \end{aligned}$$

In the new basis $x_2\mathbf{b}'_1 + y_2\mathbf{b}'_2 + z_2\mathbf{b}'_3 = (x_2 - \frac{1}{2}y_2 + z_2)\mathbf{b}_1 + y_2\mathbf{b}_2 + z_2\mathbf{b}_3$, which tells us that

$$q(x_2\mathbf{b}'_1 + y_2\mathbf{b}'_2 + z_2\mathbf{b}'_3) = x_2^2 - \frac{1}{4}y_2^2 + 4y_2z_2 - z_2^2.$$

so the matrix of q with respect to the \mathbf{b}' basis is

$$B' = \begin{pmatrix} 1 & 0 & 0 \\ 0 & -1/4 & 2 \\ 0 & 2 & -1 \end{pmatrix} \stackrel{\heartsuit}{=} Q^T B Q \stackrel{\heartsuit}{=} (P')^T A P'.$$

4 Bilinear Maps and Quadratic Forms

Third change of basis: Now we are in Step 3 of the proof, concentrating on the bottom right 2×2 block. We must change the second and third basis vectors. Any subsequent changes of basis we make will keep the first basis vector unchanged. We have

$$q(y_2 \mathbf{b}'_2 + z_2 \mathbf{b}'_3) = -\frac{1}{4}y_2^2 + 4y_2z_2 - z_2^2,$$

the “leftover terms” of the bottom right corner. This is a 2-variable quadratic form.

Since $q(\mathbf{b}'_2) = -1/4 \neq 0$, we don't need to do anything for Step 1 of the proof. Using Step 2 of the proof, we replace $\mathbf{b}'_1, \mathbf{b}'_2, \mathbf{b}'_3$ by another new basis \mathbf{b}'' :

$$\mathbf{b}''_1 := \mathbf{b}'_1, \mathbf{b}''_2 := \mathbf{b}'_2, \mathbf{b}''_3 := \mathbf{b}'_3 - \frac{2}{-1/4}\mathbf{b}'_2 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} + 8 \begin{pmatrix} -1/2 \\ 1/2 \\ 0 \end{pmatrix} = \begin{pmatrix} -3 \\ 5 \\ 1 \end{pmatrix}.$$

So the change of basis matrix from \mathbf{e} to \mathbf{b}'' is

$$P'' = \begin{pmatrix} 1 & -1/2 & -3 \\ 1 & 1/2 & 5 \\ 0 & 0 & 1 \end{pmatrix} \stackrel{\heartsuit}{=} P'Q' \text{ where } Q' = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 8 \\ 0 & 0 & 1 \end{pmatrix}.$$

This corresponds, of course, to the completing-the-square operation

$$-\frac{1}{4}y_2^2 + 4y_2z_2 - z_2^2 = -\frac{1}{4}(y_2 - 8z_2)^2 + 15z_2^2.$$

So the matrix of q is now

$$B'' = \begin{pmatrix} 1 & 0 & 0 \\ 0 & -1/4 & 0 \\ 0 & 0 & 15 \end{pmatrix} \stackrel{\heartsuit}{=} (Q')^T B' Q' \stackrel{\heartsuit}{=} (P'')^T A P''.$$

This is diagonal, so we're done: the matrix of q in the basis $\mathbf{b}''_1, \mathbf{b}''_2, \mathbf{b}''_3$ is the diagonal matrix B'' .

Notice that the choice of “good” basis, and the resulting “good” matrix, are extremely far from unique. For instance, in the example above we could have replaced \mathbf{b}''_2 with $2\mathbf{b}''_2$ to get the (perhaps nicer) matrix

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 15 \end{pmatrix}.$$

In the case $K = \mathbb{C}$, we can do even better. After reducing q to the form $q(\mathbf{v}) = \sum_{i=1}^n \alpha_{ii} x_i^2$, we can permute the coordinates if necessary to get $\alpha_{ii} \neq 0$ for $1 \leq i \leq r$ and $\alpha_{ii} = 0$ for $r+1 \leq i \leq n$, where $r = \text{rank}(q)$. We can then make a further change of coordinates $x'_i = \sqrt{\alpha_{ii}} x_i$ ($1 \leq i \leq r$), giving $q(\mathbf{v}) = \sum_{i=1}^r (x'_i)^2$. Hence we have proved:

4 Bilinear Maps and Quadratic Forms

Proposition 4.4.2. *A quadratic form q over \mathbb{C} has the form $q(\mathbf{v}) = \sum_{i=1}^r x_i^2$ with respect to a suitable basis, where $r = \text{rank}(q)$.*

Equivalently, given a symmetric matrix $A \in \mathbb{C}^{n,n}$, there is an invertible matrix $P \in \mathbb{C}^{n,n}$ such that $P^T A P = B$, where $B = (\beta_{ij})$ is a diagonal matrix with $\beta_{ii} = 1$ for $1 \leq i \leq r$, $\beta_{ii} = 0$ for $r+1 \leq i \leq n$, and $r = \text{rank}(A)$.

In particular, up to changes of basis, a quadratic form on \mathbb{C}^n is uniquely determined by its rank. We say the rank is the only *invariant* of a quadratic form over \mathbb{C} .

When $K = \mathbb{R}$, we cannot take square roots of negative numbers, but we can replace each positive α_i by 1 and each negative α_i by -1 to get:

Proposition 4.4.3 (Sylvester's Theorem). *A quadratic form q over \mathbb{R} has the form $q(\mathbf{v}) = \sum_{i=1}^t x_i^2 - \sum_{i=1}^u x_{t+i}^2$ with respect to a suitable basis, where $t + u = \text{rank}(q)$.*

Equivalently, given a symmetric matrix $A \in \mathbb{R}^{n,n}$, there is an invertible matrix $P \in \mathbb{R}^{n,n}$ such that $P^T A P = B$, where $B = (\beta_{ij})$ is a diagonal matrix with $\beta_{ii} = 1$ for $1 \leq i \leq t$, $\beta_{ii} = -1$ for $t+1 \leq i \leq t+u$, and $\beta_{ii} = 0$ for $t+u+1 \leq i \leq n$, and $t+u = \text{rank}(A)$.

We shall now prove that the numbers t and u of positive and negative terms are invariants of q . The pair of integers (t, u) is called the *signature* of q .

Theorem 4.4.4 (Sylvester's Law of Inertia). *Suppose that q is a quadratic form on the vector space V over \mathbb{R} , and that $\mathbf{e}_1, \dots, \mathbf{e}_n$ and $\mathbf{e}'_1, \dots, \mathbf{e}'_n$ are two bases of V such that*

$$q(x_1 \mathbf{e}_1 + \dots + x_n \mathbf{e}_n) = \sum_{i=1}^t x_i^2 - \sum_{i=1}^u x_{t+i}^2$$

and

$$q(x_1 \mathbf{e}'_1 + \dots + x_n \mathbf{e}'_n) = \sum_{i=1}^{t'} x_i^2 - \sum_{i=1}^{u'} x_{t'+i}^2.$$

Then $t = t'$ and $u = u'$.

Proof. We know that $t + u = t' + u' = \text{rank}(q)$, so it is enough to prove that $t = t'$. Suppose not; by symmetry we may suppose that $t > t'$.

Let V_1 be the span of $\mathbf{e}_1, \dots, \mathbf{e}_t$, and let V_2 be the span of $\mathbf{e}'_{t'+1}, \dots, \mathbf{e}'_n$. Then for any non-zero $\mathbf{v} \in V_1$ we have $q(\mathbf{v}) > 0$; while for any $\mathbf{w} \in V_2$ we have $q(\mathbf{w}) \leq 0$. So there cannot be any non-zero $\mathbf{v} \in V_1 \cap V_2$.

On the other hand, we have $\dim(V_1) = t$ and $\dim(V_2) = n - t'$. It was proved in MA106 that

$$\dim(V_1 + V_2) = \dim(V_1) + \dim(V_2) - \dim(V_1 \cap V_2),$$

so

$$\dim(V_1 \cap V_2) = t + (n - t') - \dim(V_1 + V_2) = (t - t') + n - \dim(V_1 + V_2) > 0.$$

The last inequality follows from our assumption on $t - t'$ and the fact $V_1 + V_2$ is a subspace of V and thus has dimension at most n . Since we have shown that $V_1 \cap V_2 = \{0\}$, this is a contradiction, which completes the proof. \square

4 Bilinear Maps and Quadratic Forms

Remark. Notice that any non-zero $x \in \mathbb{R}$ is either equal to a square, or -1 times a square, but not both. This property is shared by the finite field \mathbb{F}_7 of integers mod 7, so any quadratic form over \mathbb{F}_7 can be written as a diagonal matrix with only 0's, 1's and -1 's down the diagonal (i.e. Sylvester's Theorem holds over \mathbb{F}_7). But Sylvester's law of inertia isn't valid in \mathbb{F}_7 : in fact, we have

$$\begin{pmatrix} 2 & 3 \\ 4 & 2 \end{pmatrix}^T \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 2 & 3 \\ 4 & 2 \end{pmatrix} = \begin{pmatrix} 20 & 14 \\ 14 & 20 \end{pmatrix} = \begin{pmatrix} -1 & 0 \\ 0 & -1 \end{pmatrix},$$

so the same form has signature $(2, 0)$ and $(0, 2)$! The proof breaks down because there's no good notion of a "positive" element of \mathbb{F}_7 , so a sum of non-zero squares can be zero (the easiest example is $1^2 + 2^2 + 3^2 = 0$). So Sylvester's law of inertia is really using something quite special about \mathbb{R} .

4.5 Euclidean spaces, orthonormal bases and the Gram–Schmidt process

In this section, we're going to suppose $K = \mathbb{R}$. As usual, we let V be an n -dimensional vector space over K , and we let q be a quadratic form on V , with associated symmetric bilinear form τ .

Definition 4.5.1. The quadratic form q is said to be *positive definite* if $q(\mathbf{v}) > 0$ for all $0 \neq \mathbf{v} \in V$.

It is clear that this is the case if and only if $t = n$ and $u = 0$ in Proposition 4.4.3; that is, if q has signature $(n, 0)$.

The associated symmetric bilinear form τ is also called positive definite when q is.

Definition 4.5.2. A vector space V over \mathbb{R} together with a positive definite symmetric bilinear form τ is called a *euclidean space*.

In this case, Proposition 4.4.3 says that there is a basis $\{\mathbf{e}_i\}$ of V with respect to which $\tau(\mathbf{e}_i, \mathbf{e}_j) = \delta_{ij}$, where

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j. \end{cases}$$

(so the matrix A of q is the identity matrix I_n .) We call a basis of a euclidean space V with this property an *orthonormal basis* of V .

(More generally, any set $\mathbf{v}_1, \dots, \mathbf{v}_r$ of vectors in V , not necessarily a basis, will be said to be *orthonormal* if $\tau(\mathbf{v}_i, \mathbf{v}_j) = \delta_{ij}$ for $1 \leq i, j \leq r$.)

We shall assume from now on that V is a euclidean space, and that we have chosen an orthonormal basis $\mathbf{e}_1, \dots, \mathbf{e}_n$. Then τ corresponds to the standard dot product and we shall write $\mathbf{v} \cdot \mathbf{w}$ instead of $\tau(\mathbf{v}, \mathbf{w})$.

Note that $\mathbf{v} \cdot \mathbf{w} = \underline{\mathbf{v}}^T \underline{\mathbf{w}}$ where, as usual, $\underline{\mathbf{v}}$ and $\underline{\mathbf{w}}$ are the column vectors associated with \mathbf{v} and \mathbf{w} .

For $\mathbf{v} \in V$, define $|\mathbf{v}| = \sqrt{\mathbf{v} \cdot \mathbf{v}}$. Then $|\mathbf{v}|$ is the length of \mathbf{v} . Hence the length, and also the cosine $\mathbf{v} \cdot \mathbf{w} / (|\mathbf{v}| |\mathbf{w}|)$ of the angle between two vectors can be defined in terms of the scalar product.

4 Bilinear Maps and Quadratic Forms

Thus a set of vectors is orthonormal if the vectors all have length 1 and are at right angles to each other.

The following theorem tells us that we can always complete a set of orthonormal vectors to an orthonormal basis.

Theorem 4.5.3 (Gram-Schmidt process). *Let V be a euclidean space of dimension n , and suppose that, for some r with $0 \leq r \leq n$, $\mathbf{f}_1, \dots, \mathbf{f}_r$ are vectors in V such that*

$$\mathbf{f}_i \cdot \mathbf{f}_j = \delta_{ij} \quad \text{for } 1 \leq i, j \leq r. \quad (*)$$

Then $\mathbf{f}_1, \dots, \mathbf{f}_r$ can be extended to an orthonormal basis $\mathbf{f}_1, \dots, \mathbf{f}_n$ of V .

Proof. We prove first that $\mathbf{f}_1, \dots, \mathbf{f}_r$ are linearly independent. Suppose that $\sum_{i=1}^r \alpha_i \mathbf{f}_i = \mathbf{0}$ for some $\alpha_1, \dots, \alpha_r \in \mathbb{R}$. Then, for each j with $1 \leq j \leq r$, the scalar product of the left hand side of this equation with \mathbf{f}_j is $\sum_{i=1}^r \alpha_i \mathbf{f}_j \cdot \mathbf{f}_i = \alpha_j$, by (*). Since $\mathbf{f}_j \cdot \mathbf{0} = 0$, this implies that $\alpha_j = 0$ for all j , so the \mathbf{f}_i are linearly independent.

The proof of the theorem will be by induction on $n - r$. We can start the induction with the case $n - r = 0$, when $r = n$, and there is nothing to prove. So assume that $n - r > 0$; i.e. that $r < n$. By a result from MA106, we can extend any linearly independent set of vectors to a basis of V , so there is a basis $\mathbf{f}_1, \dots, \mathbf{f}_r, \mathbf{g}_{r+1}, \dots, \mathbf{g}_n$ of V containing the \mathbf{f}_i . The trick is to define

$$\mathbf{f}'_{r+1} = \mathbf{g}_{r+1} - \sum_{i=1}^r (\mathbf{f}_i \cdot \mathbf{g}_{r+1}) \mathbf{f}_i.$$

If we take the scalar product of this equation with \mathbf{f}_j for some $1 \leq j \leq r$, then we get

$$\mathbf{f}_j \cdot \mathbf{f}'_{r+1} = \mathbf{f}_j \cdot \mathbf{g}_{r+1} - \sum_{i=1}^r (\mathbf{f}_i \cdot \mathbf{g}_{r+1}) (\mathbf{f}_j \cdot \mathbf{f}_i)$$

and then, by (*), $\mathbf{f}_j \cdot \mathbf{f}_i$ is non-zero only when $j = i$, so the sum on the right hand side simplifies to $\mathbf{f}_j \cdot \mathbf{g}_{r+1}$, and the whole equation simplifies to

$$\mathbf{f}_j \cdot \mathbf{f}'_{r+1} = \mathbf{f}_j \cdot \mathbf{g}_{r+1} - \mathbf{f}_j \cdot \mathbf{g}_{r+1} = 0.$$

The vector \mathbf{f}'_{r+1} is non-zero by linear independence of the basis, and if we define $\mathbf{f}_{r+1} = \mathbf{f}'_{r+1} / |\mathbf{f}'_{r+1}|$, then we still have $\mathbf{f}_j \cdot \mathbf{f}_{r+1} = 0$ for $1 \leq j \leq r$, and we also have $\mathbf{f}_{r+1} \cdot \mathbf{f}_{r+1} = 1$. Hence $\mathbf{f}_1, \dots, \mathbf{f}_{r+1}$ satisfy the equations (*), and the result follows by invoking our inductive hypothesis. \square

Note that this proof is constructive. In fact it shows us that given any basis of a euclidean space we can 'correct it' to an orthonormal basis, as in the following example.

Example. Let $V = \mathbb{R}^3$ with the standard dot product. It is straightforward to check that

$\begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix}$ is a basis for V but it is not orthonormal. Let's use the Gram-Schmidt

process to fix that by taking $r = 0$ and $g_1 = \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}, g_2 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$ and $g_3 = \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix}$.

4 Bilinear Maps and Quadratic Forms

Then $f'_1 := g_1$ and so $f_1 = f'_1/|f'_1| = \frac{1}{\sqrt{3}} \begin{pmatrix} 1 \\ -1 \\ 1 \end{pmatrix}$,

$f'_2 := g_2 - (f_1 \cdot g_2)f_1 = g_2 - \frac{2}{\sqrt{3}}f_1 = \frac{1}{3} \begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix}$ and so $f_2 = \frac{1}{\sqrt{6}} \begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix}$,

$f'_3 := g_3 - (f_1 \cdot g_3)f_1 - (f_2 \cdot g_3)f_2 = g_3 - \frac{2}{\sqrt{3}}f_1 - \frac{5}{\sqrt{6}}f_2 = \frac{1}{2} \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}$ and so $f_3 = \frac{1}{\sqrt{2}} \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}$.

thus we have now got an orthonormal basis f_1, f_2, f_3 (always good to check this at the end!).

4.6 Orthogonal transformations

If we're working with a euclidean space V , we know that there is a sensible definition of length of a vector and the angle between vectors is; so we might want to consider transformations from V to itself that preserve lengths and angles – they play nicely with the geometry of the space.

Definition 4.6.1. A linear map $T:V \rightarrow V$ is said to be *orthogonal* if it preserves the scalar product on V . That is, if $T(\mathbf{v}) \cdot T(\mathbf{w}) = \mathbf{v} \cdot \mathbf{w}$ for all $\mathbf{v}, \mathbf{w} \in V$.

Since length and angle can be defined in terms of the scalar product, an orthogonal linear map preserves distance and angle, so geometrically it is a rigid map. In \mathbb{R}^2 , for example, an orthogonal map is either a rotation about the origin, or a reflection about a line through the origin.

If A is the matrix of T (with respect to some orthonormal basis), then $T(\mathbf{v}) = A\mathbf{v}$ and so

$$T(\mathbf{v}) \cdot T(\mathbf{w}) = \mathbf{v}^T A^T A \mathbf{w}.$$

Hence T is orthogonal (the right hand side equals $\mathbf{v} \cdot \mathbf{w}$) if and only if $A^T A = I_n$, or equivalently if $A^T = A^{-1}$.

Definition 4.6.2. An $n \times n$ matrix is called *orthogonal* if $A^T A = I_n$.

So we have proved:

Proposition 4.6.3. A linear map $T : V \rightarrow V$ is orthogonal if and only if its matrix A (with respect to an orthonormal basis of V) is orthogonal.

Incidentally, the fact that $A^T A = I_n$ tells us that A (and hence T) is invertible, so $\det(A)$ is non-zero. In fact we can do a little better than that:

Proposition 4.6.4. An orthogonal matrix has determinant ± 1 .

4 Bilinear Maps and Quadratic Forms

Proof. We have $A^T A = I_n$, so $\det(A^T A) = \det(I_n) = 1$.

On the other hand, $\det(A^T A) = \det(A^T) \det(A) = (\det A)^2$. So $(\det A)^2 = 1$, implying that $\det A = \pm 1$. \square

Example. For any $\theta \in \mathbb{R}$, let $A = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$. (This represents a anticlockwise rotation through an angle θ .) Then it is easily checked that $A^T A = A A^T = I_2$.

One can check that every orthogonal 2×2 matrix with determinant $+1$ is a rotation by some angle θ , and similarly that any orthogonal 2×2 matrix of $\det -1$ is a reflection in some line through the origin. In higher dimensions the taxonomy of orthogonal matrices is a bit more complicated – we’ll revisit this in a later section of the course.

Notice that the columns of A are mutually orthogonal vectors of length 1, and the same applies to the rows of A . Let $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n$ be the columns of the matrix A . As we observed in §1, \mathbf{c}_i is equal to the column vector representing $T(\mathbf{e}_i)$. In other words, if $T(\mathbf{e}_i) = \mathbf{f}_i$, say, then $\underline{\mathbf{f}}_i = \mathbf{c}_i$.

Since the (i, j) -th entry of $A^T A$ is $\mathbf{c}_i^T \mathbf{c}_j = \mathbf{f}_i \cdot \mathbf{f}_j$, we see that T and A are orthogonal if and only if

$$\mathbf{f}_i \cdot \mathbf{f}_i = 1 \text{ and } \mathbf{f}_i \cdot \mathbf{f}_j = 0 \ (i \neq j), \quad 1 \leq i, j \leq n. \quad (*)$$

By Proposition 4.6.4, an orthogonal linear map is invertible, so $T(\mathbf{e}_i)$ ($1 \leq i \leq n$) forms a basis of V , and we have:

Proposition 4.6.5. *A linear map T is orthogonal if and only if $T(\mathbf{e}_1), \dots, T(\mathbf{e}_n)$ is an orthonormal basis of V .*

Here’s a pretty application of the Gram-Schmidt process and orthogonal matrices. Notice that our proof of Gram-Schmidt actually proved a little more: we showed that if $\mathbf{f}_1, \dots, \mathbf{f}_r$ is an orthonormal set, and $\mathbf{g}_{r+1}, \dots, \mathbf{g}_n$ is any way of completing the \mathbf{f} ’s to a basis of V , then we can find $\mathbf{f}_{r+1}, \dots, \mathbf{f}_n$ such that

- $\mathbf{f}_1, \dots, \mathbf{f}_n$ is an orthonormal basis,
- for each $r + 1 \leq i \leq n$, \mathbf{f}_i is in the linear span of $\mathbf{f}_1, \dots, \mathbf{f}_r, \mathbf{g}_{r+1}, \dots, \mathbf{g}_i$.

That is, we’ve arranged that the basis change matrix from $\mathbf{f}_1, \dots, \mathbf{f}_r, \mathbf{g}_{r+1}, \dots, \mathbf{g}_n$ to $\mathbf{f}_1, \dots, \mathbf{f}_n$ looks like

$$\left(\begin{array}{c|c} I_r & A \\ \hline 0 & B \end{array} \right)$$

with B upper-triangular. This is most useful when $r = 0$, when it says that any basis may be modified by an upper-triangular matrix to make it orthonormal.

This slight modification of Gram-Schmidt has a very nice interpretation in terms of matrices:

Proposition 4.6.6 (QR decomposition). *Let A be any $n \times n$ real matrix. Then we can write $A = QR$ where Q is orthogonal and R is upper-triangular.*

4 Bilinear Maps and Quadratic Forms

Proof. We'll only include a proof for the case when A is invertible.

Let $\mathbf{g}_1, \dots, \mathbf{g}_n$ be the columns of A , regarded as vectors in \mathbb{R}^n . Then since A is invertible, $\mathbf{g}_1, \dots, \mathbf{g}_n$ is a basis of \mathbb{R}^n . We apply the Gram-Schmidt process to construct an orthonormal basis $\mathbf{f}_1, \dots, \mathbf{f}_n$ such that \mathbf{f}_i is in the linear span of $\mathbf{g}_1, \dots, \mathbf{g}_i$ for each i .

Let Q be the matrix whose columns are $\mathbf{f}_1, \dots, \mathbf{f}_n$. Then Q is an orthogonal matrix, since its columns are orthonormal vectors; and there are real numbers r_{ij} such that

$$\begin{aligned}\mathbf{g}_1 &= r_{11}\mathbf{f}_1 \\ \mathbf{g}_2 &= r_{12}\mathbf{f}_1 + r_{22}\mathbf{f}_2 \\ \mathbf{g}_3 &= r_{13}\mathbf{f}_1 + r_{23}\mathbf{f}_2 + r_{33}\mathbf{f}_3 \\ &\vdots\end{aligned}$$

In other words we have

$$A = QR$$

where R is the upper-triangular matrix with entries r_{ij} . □

To deal with the case when A isn't invertible (so the columns of A no longer form a basis) we can do the following. Show that any matrix A can be written as $A = BR'$ where B is invertible and R' is upper triangular; then writing $B = QR$ we have $A = QRR'$, and RR' is also upper-triangular.

Example. Consider the matrix

$$A = \begin{pmatrix} -1 & 0 & -2 \\ 2 & 0 & -1 \\ 0 & -2 & -2 \end{pmatrix}.$$

We have $\det(A) = 10$, so A is non-singular. Let $\mathbf{g}_1, \mathbf{g}_2, \mathbf{g}_3$ be the columns of A .

Then $|\mathbf{g}_1| = \sqrt{5}$, so

$$\mathbf{f}_1 = \frac{\mathbf{g}_1}{|\mathbf{g}_1|} = \begin{pmatrix} -1/\sqrt{5} \\ 2/\sqrt{5} \\ 0 \end{pmatrix}.$$

For the next step, we take $\mathbf{f}'_2 = \mathbf{g}_2 - (\mathbf{f}_1 \cdot \mathbf{g}_2)\mathbf{f}_1 = \mathbf{g}_2$, since $\mathbf{f}_1 \cdot \mathbf{g}_2 = 0$. So

$$\mathbf{f}_2 = \frac{\mathbf{g}_2}{|\mathbf{g}_2|} = \begin{pmatrix} 0 \\ 0 \\ -1 \end{pmatrix}.$$

For the final step, we take the vector

$$\mathbf{f}'_3 = \mathbf{g}_3 - (\mathbf{f}_1 \cdot \mathbf{g}_3)\mathbf{f}_1 - (\mathbf{f}_2 \cdot \mathbf{g}_3)\mathbf{f}_2.$$

We have

$$\mathbf{f}_1 \cdot \mathbf{g}_3 = \begin{pmatrix} -1/\sqrt{5} \\ 2/\sqrt{5} \\ 0 \end{pmatrix} \cdot \begin{pmatrix} -2 \\ -1 \\ -2 \end{pmatrix} = 0, \quad \mathbf{f}_2 \cdot \mathbf{g}_3 = \begin{pmatrix} 0 \\ 0 \\ -1 \end{pmatrix} \cdot \begin{pmatrix} -2 \\ -1 \\ -2 \end{pmatrix} = 2.$$

4 Bilinear Maps and Quadratic Forms

So $\mathbf{f}'_3 = \mathbf{g}_3 - 2\mathbf{f}_2 = \begin{pmatrix} -2 \\ -1 \\ 0 \end{pmatrix}$. We have $|\mathbf{f}'_3| = \sqrt{5}$ again, so

$$\mathbf{f}_3 = \frac{\mathbf{f}'_3}{\sqrt{5}} = \begin{pmatrix} -2/\sqrt{5} \\ -1/\sqrt{5} \\ 0 \end{pmatrix}.$$

Thus Q is the matrix whose columns are $\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3$, that is

$$Q = \begin{pmatrix} -1/\sqrt{5} & 0 & -2/\sqrt{5} \\ 2/\sqrt{5} & 0 & -1/\sqrt{5} \\ 0 & -1 & 0 \end{pmatrix}.$$

and we have

$$\mathbf{g}_1 = \sqrt{5}\mathbf{f}_1, \quad \mathbf{g}_2 = 2\mathbf{f}_2, \quad \mathbf{g}_3 = 2\mathbf{f}_2 + \sqrt{5}\mathbf{f}_3$$

so $A = QR$ where

$$R = \begin{pmatrix} \sqrt{5} & 0 & 0 \\ 0 & 2 & 2 \\ 0 & 0 & \sqrt{5} \end{pmatrix}.$$

The QR decomposition theorem is a very important technique in numerical calculations with matrices. For instance, QR decomposition gives a quick way of inverting matrices. If $A = QR$, then $A^{-1} = R^{-1}Q^{-1}$. Inverting orthogonal matrices is trivial, as the inverse is just the transpose; inverting upper-triangular matrices is also pretty easy, so we can compute the inverse of A this way, without having to compute the determinant.

4.7 Nice orthonormal bases

Suppose we have a euclidean space V , and a linear operator $T : V \rightarrow V$ (in case it is not clear linear operator just means a linear map from V to V - it is an operator on V), or a quadratic form q on V (not necessarily the same as the one giving V its euclidean structure). Can we always find an orthonormal basis of V making the matrix of q look reasonably nice? Notice that we're juggling two quadratic forms here - we're trying to make the matrix of q look nice while simultaneously keeping the matrix of the original quadratic form as the identity.

Oddly enough, this is *also* a question about linear operators. Given any bilinear form τ on V (not necessarily symmetric), there's a uniquely determined linear operator T on V such that

$$\tau(\mathbf{v}, \mathbf{w}) = \mathbf{v} \cdot T(\mathbf{w}).$$

Note that T is just the linear operator corresponding to the matrix A of τ , with i, j th entry $\tau(e_i, e_j)$ for the standard basis e_1, \dots, e_n of V . Indeed, then $\tau(\mathbf{v}, \mathbf{w}) = \underline{\mathbf{v}}^T A \underline{\mathbf{w}} = \mathbf{v} \cdot T(\mathbf{w})$ (this will be true if we used any orthonormal basis of V but not in general).

4 Bilinear Maps and Quadratic Forms

Conversely, any linear operator T determines a bilinear form τ by the same formula. Indeed, the bilinearity follows from the bilinearity of \cdot and the linearity of T (check this!). So once we've fixed a "starting" bilinear form (the positive definite symmetric bilinear form), we can get any other bilinear form τ on $V \times V$ from this via a linear operator, and this gives us a bijection between bilinear forms and linear operators. Moreover, the matrix of T , in an orthonormal basis of V , is just the matrix of τ . When we change basis by an orthogonal matrix (to get a new orthonormal basis), the matrix of T changes by $A \mapsto P^{-1}AP$, and the matrix of τ changes by $A \mapsto P^TAP$, but this is OK since $P^T = P^{-1}$ for orthogonal matrices!

In particular, if T is any linear operator, then $(\mathbf{v}, \mathbf{w}) \mapsto (T\mathbf{v}) \cdot \mathbf{w}$ is certainly a bilinear form; so there must be some linear operator S such that

$$(T\mathbf{v}) \cdot \mathbf{w} = \mathbf{v} \cdot (S\mathbf{w}) \tag{*}$$

for all \mathbf{v} and \mathbf{w} .

Definition 4.7.1. If $T : V \rightarrow V$ is a linear operator on a euclidean space V , then the unique linear map S such that (*) holds is called the *adjoint* of T . We write this as T^* .

If we have chosen an orthonormal basis, then the matrix of T^* is just the transpose of the matrix of T . It follows from this that a linear operator is orthogonal if and only if $T^* = T^{-1}$; one can also prove this directly from the definition.

We say T is *selfadjoint* if $T^* = T$, or equivalently if the bilinear form $\tau(\mathbf{v}, \mathbf{w}) = \mathbf{v} \cdot (T\mathbf{w})$ is symmetric. Notice that 'selfadjointness', like 'orthogonalness', is something that only makes sense for linear operators on euclidean spaces; it doesn't make sense to ask if a linear operator on a general vector space is selfadjoint. It should be clear that T is selfadjoint if and only if its matrix in an orthonormal basis of V is a symmetric matrix.

So if V is a euclidean space of dimension n , the following problems are all actually the same:

- given a quadratic form q on V , find an orthonormal basis of V making the matrix of q as nice as possible;
- given a selfadjoint linear operator T on V , find an orthonormal basis of V making the matrix of T as nice as possible;
- given an $n \times n$ symmetric real matrix A , find an orthogonal matrix P such that P^TAP is as nice as possible.

First, we'll warm up by proving a proposition which we'll need in proving the main result solving these equivalent problems.

Proposition 4.7.2. *Let A be an $n \times n$ real symmetric matrix. Then A has an eigenvalue in \mathbb{R} , and all complex eigenvalues of A lie in \mathbb{R} .*

Proof. (To simplify the notation, we will write just \mathbf{v} for a column vector $\underline{\mathbf{v}}$ in this proof.)

The characteristic equation $\det(A - xI_n) = 0$ is a polynomial equation of degree n in x , and since \mathbb{C} is an algebraically closed field, it certainly has a root $\lambda \in \mathbb{C}$, which is an eigenvalue for

4 Bilinear Maps and Quadratic Forms

A if we regard A as a matrix over \mathbb{C} . We shall prove that any such λ lies in \mathbb{R} , which will prove the proposition.

For a column vector \mathbf{v} or matrix B over \mathbb{C} , we denote by $\bar{\mathbf{v}}$ or \bar{B} the result of replacing all entries of \mathbf{v} or B by their complex conjugates. Since the entries of A lie in \mathbb{R} , we have $\bar{A} = A$.

Let \mathbf{v} be a complex eigenvector associated with λ . Then

$$A\mathbf{v} = \lambda\mathbf{v} \tag{1}$$

so, taking complex conjugates and using $\bar{A} = A$, we get

$$A\bar{\mathbf{v}} = \bar{\lambda}\bar{\mathbf{v}}. \tag{2}$$

Transposing (1) and using $A^T = A$ gives

$$\mathbf{v}^T A = \lambda\mathbf{v}^T, \tag{3}$$

so by (2) and (3) we have

$$\lambda\mathbf{v}^T\bar{\mathbf{v}} = \mathbf{v}^T A\bar{\mathbf{v}} = \bar{\lambda}\mathbf{v}^T\bar{\mathbf{v}}.$$

But if $\mathbf{v} = (\alpha_1, \alpha_2, \dots, \alpha_n)^T$, then $\mathbf{v}^T\bar{\mathbf{v}} = \alpha_1\bar{\alpha}_1 + \dots + \alpha_n\bar{\alpha}_n$, which is a non-zero real number (eigenvectors are non-zero by definition). Thus $\lambda = \bar{\lambda}$, so $\lambda \in \mathbb{R}$. □

Now let's prove the main theorem of this section.

Theorem 4.7.3. *Let V be a euclidean space of dimension n . Then:*

- *Given any quadratic form q on V , there is an orthonormal basis $\mathbf{f}_1, \dots, \mathbf{f}_n$ of V and constants $\alpha_1, \dots, \alpha_n$, uniquely determined up to reordering, such that*

$$q(x_1\mathbf{f}_1 + \dots + x_n\mathbf{f}_n) = \sum_{i=1}^n \alpha_i(x_i)^2$$

for all $x_1, \dots, x_n \in \mathbb{R}$.

- *Given any linear operator $T : V \rightarrow V$ which is selfadjoint, there is an orthonormal basis $\mathbf{f}_1, \dots, \mathbf{f}_n$ of V consisting of eigenvectors of T .*
- *Given any $n \times n$ real symmetric matrix A , there is an orthogonal matrix P such that $P^T A P = P^{-1} A P$ is a diagonal matrix.*

Proof. We've already seen that these three statements are equivalent to each other, so we can prove whichever one of them we like. Notice that in the second and third forms of the statement, it's clear that the diagonal matrix we obtain is similar to the original one; that tells us that in the first statement the constants $\alpha_1, \dots, \alpha_n$ are uniquely determined (possibly up to re-ordering).

We'll prove the second statement using induction on $n = \dim V$. If $n = 0$ there is nothing to prove, so let's assume the proposition holds for $n - 1$.

4 Bilinear Maps and Quadratic Forms

Let T be our linear operator. By Proposition 4.7.2, T has an eigenvalue in \mathbb{R} . Let \mathbf{v} be a corresponding eigenvector in V . Then $\mathbf{f}_1 = \mathbf{v}/|\mathbf{v}|$ is also an eigenvector, and $|\mathbf{f}_1| = 1$. Let α_1 be the corresponding eigenvalue.

We consider the space $W = \{\mathbf{w} \in V : \mathbf{w} \cdot \mathbf{f}_1 = 0\}$. Since W is the kernel of a surjective linear map

$$V \longrightarrow \mathbb{R}, \quad \mathbf{v} \mapsto \mathbf{v} \cdot \mathbf{f}_1,$$

it is a subspace of V of dimension $n - 1$. We claim that T maps W into itself. So suppose $\mathbf{w} \in W$; we want to show that $T(\mathbf{w}) \in W$ also.

We have

$$T(\mathbf{w}) \cdot \mathbf{f}_1 = \mathbf{w} \cdot T(\mathbf{f}_1)$$

since T is selfadjoint. But we know that $T(\mathbf{f}_1) = \alpha_1 \mathbf{f}_1$, so it follows that

$$T(\mathbf{w}) \cdot \mathbf{f}_1 = \alpha_1 (\mathbf{w} \cdot \mathbf{f}_1) = 0,$$

since $\mathbf{w} \in W$ so $\mathbf{w} \cdot \mathbf{f}_1 = 0$ since $\alpha_1 \neq 0$.

So T maps W into itself. Moreover, W is a euclidean space of dimension $n - 1$, so we may apply the induction hypothesis to the restriction of T to W . This gives us an orthonormal basis $\mathbf{f}_2, \dots, \mathbf{f}_n$ of W consisting of eigenvectors of T . By definition of W , \mathbf{f}_1 is orthogonal to $\mathbf{f}_2, \dots, \mathbf{f}_n$ and it follows that $\mathbf{f}_1, \dots, \mathbf{f}_n$ is an orthonormal basis of V , consisting of eigenvectors of T . \square

Although it is not used in the proof of the theorem above, the following proposition is useful when calculating examples. It helps us to write down more vectors in the final orthonormal basis immediately, without having to use Theorem 4.5.3 repeatedly.

Proposition 4.7.4. *Let A be a real symmetric matrix, and let λ_1, λ_2 be two distinct eigenvalues of A , with corresponding eigenvectors $\mathbf{v}_1, \mathbf{v}_2$. Then $\mathbf{v}_1 \cdot \mathbf{v}_2 = 0$.*

Proof. (As in Proposition 4.7.2, we will write \mathbf{v} rather than $\underline{\mathbf{v}}$ for a column vector in this proof. So $\mathbf{v}_1 \cdot \mathbf{v}_2$ is the same as $\mathbf{v}_1^T \mathbf{v}_2$.) We have

$$A\mathbf{v}_1 = \lambda_1 \mathbf{v}_1, \tag{1}$$

$$A\mathbf{v}_2 = \lambda_2 \mathbf{v}_2. \tag{2}$$

The trick is now to look at the expression $\mathbf{v}_1^T A\mathbf{v}_2$. On the one hand, by (2) we have

$$\mathbf{v}_1^T A\mathbf{v}_2 = \mathbf{v}_1 \cdot (A\mathbf{v}_2) = \mathbf{v}_1^T (\lambda_2 \mathbf{v}_2) = \lambda_2 (\mathbf{v}_1 \cdot \mathbf{v}_2). \tag{3}$$

On the other hand, $A^T = A$, so $\mathbf{v}_1^T A = \mathbf{v}_1^T A^T = (A\mathbf{v}_1)^T$, so using (1) we have

$$\mathbf{v}_1^T A\mathbf{v}_2 = (A\mathbf{v}_1)^T \mathbf{v}_2 = (\lambda_1 \mathbf{v}_1^T) \mathbf{v}_2 = \lambda_1 (\mathbf{v}_1 \cdot \mathbf{v}_2). \tag{4}$$

Comparing (3) and (4), we have $(\lambda_2 - \lambda_1)(\mathbf{v}_1 \cdot \mathbf{v}_2) = 0$. Since $\lambda_2 - \lambda_1 \neq 0$ by assumption, we have $\mathbf{v}_1^T \mathbf{v}_2 = 0$. \square

4 Bilinear Maps and Quadratic Forms

Example 16. Let $n = 2$ and let A be the symmetric matrix $A = \begin{pmatrix} 1 & 3 \\ 3 & 1 \end{pmatrix}$. Then

$$\det(A - xI_2) = (1 - x)^2 - 9 = x^2 - 2x - 8 = (x - 4)(x + 2),$$

so the eigenvalues of A are 4 and -2 . Solving $A\mathbf{v} = \lambda\mathbf{v}$ for $\lambda = 4$ and -2 , we find corresponding eigenvectors $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$ and $\begin{pmatrix} 1 \\ -1 \end{pmatrix}$. Proposition 4.7.4 tells us that these vectors are orthogonal to each other (which we can of course check directly!), so if we divide them by their lengths to give vectors of length 1, giving $\begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}$ and $\begin{pmatrix} \frac{1}{\sqrt{2}} \\ \frac{-1}{\sqrt{2}} \end{pmatrix}$ then we get an orthonormal basis consisting of eigenvectors of A , which is what we want. The corresponding basis change matrix P has these vectors as columns, so $P = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \end{pmatrix}$, and we can check that $P^T P = I_2$ (i.e. P is orthogonal) and that $P^T A P = \begin{pmatrix} 4 & 0 \\ 0 & -2 \end{pmatrix}$.

Example 17. Let's do an example of the "quadratic form" version of the above theorem. Let $n = 3$ and

$$q(\mathbf{v}) = 3x^2 + 6y^2 + 3z^2 - 4xy - 4yz + 2xz,$$

$$\text{so } A = \begin{pmatrix} 3 & -2 & 1 \\ -2 & 6 & -2 \\ 1 & -2 & 3 \end{pmatrix}.$$

Then, expanding by the first row,

$$\begin{aligned} \det(A - xI_3) &= (3 - x)(6 - x)(3 - x) - 4(3 - x) - 4(3 - x) + 4 + 4 - (6 - x) \\ &= -x^3 + 12x^2 - 36x + 32 = (2 - x)(x - 8)(x - 2), \end{aligned}$$

so the eigenvalues are 2 (repeated) and 8. For the eigenvalue 8, if we solve $A\mathbf{v} = 8\mathbf{v}$ then we find a solution $\mathbf{v} = \begin{pmatrix} 1 \\ -2 \\ 1 \end{pmatrix}$. Since 2 is a repeated eigenvalue, we need two corresponding eigenvectors, which must be orthogonal to each other. The equations $A\mathbf{v} = 2\mathbf{v}$ all reduce to $a - 2b + c = 0$, and so any vector $\begin{pmatrix} a \\ b \\ c \end{pmatrix}$ satisfying this equation is an eigenvector for $\lambda = 2$. By

Proposition 4.7.4 these eigenvectors will all be orthogonal to the eigenvector for $\lambda = 8$, but we will have to choose them orthogonal to each other. We can choose the first one arbitrarily, so let's choose $\begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}$. We now need another solution that is orthogonal to this. In other words, we want a, b and c not all zero satisfying $a - 2b + c = 0$ and $a - c = 0$, and $a = b = c = 1$

4 Bilinear Maps and Quadratic Forms

is a solution. So we now have a basis $\begin{pmatrix} 1 \\ -2 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$ of three mutually orthogonal eigenvectors. To get an orthonormal basis, we just need to divide by their lengths, which are, respectively, $\sqrt{6}$, $\sqrt{2}$, and $\sqrt{3}$, and then the basis change matrix P has these vectors as columns, so

$$P = \begin{pmatrix} 1/\sqrt{6} & 1/\sqrt{2} & 1/\sqrt{3} \\ -2/\sqrt{6} & 0 & 1/\sqrt{3} \\ 1/\sqrt{6} & -1/\sqrt{2} & 1/\sqrt{3} \end{pmatrix}.$$

It can then be checked that $P^T P = I_3$ and that $P^T A P$ is the diagonal matrix with entries 8, 2, 2. So if $\mathbf{f}_1, \mathbf{f}_2, \mathbf{f}_3$ is this basis, we have

$$q(x\mathbf{f}_1 + y\mathbf{f}_2 + z\mathbf{f}_3) = 8x^2 + 2y^2 + 2z^2.$$

4.8 Applications of quadratic forms to geometry

4.8.1 Reduction of the general second degree equation

The general equation of a second degree polynomial in n variables x_1, \dots, x_n is

$$\sum_{i=1}^n \alpha_i x_i^2 + \sum_{i=1}^n \sum_{j=1}^{i-1} \alpha_{ij} x_i x_j + \sum_{i=1}^n \beta_i x_i + \gamma = 0. \quad (\dagger)$$

For fixed values of the α 's, β 's and γ , this defines a *quadric* curve or surface in n -dimensional euclidean space. To study the possible shapes of the curves and surfaces thus defined, we first simplify this equation by applying coordinate changes resulting from isometries (rigid motions) of \mathbb{R}^n ; that is, transformations that preserve distance and angle.

By Theorem 4.7.3, we can apply an orthogonal basis change (that is, an isometry of \mathbb{R}^n that fixes the origin) which has the effect of eliminating the terms $\alpha_{ij} x_i x_j$ in the above sum. To carry out this step we consider the

$$\sum_{i=1}^n \alpha_i x_i^2 + \sum_{i=1}^n \sum_{j=1}^{i-1} \alpha_{ij} x_i x_j$$

term and complete the square (we are not very interested in tracking the basis in this section). When making the change of coordinates we do then have to consider its impact on the terms in $\sum_{i=1}^n \beta_i x_i$.

For example, suppose we have $x^2 + xy + y^2 + x = 0$. Then completing the square on the $x^2 + xy + y^2$ part leads us to write the equation as $(x + \frac{1}{2}y)^2 + \frac{3}{4}y^2 + x = 0$. We then change coordinates to $x_1 = x + \frac{1}{2}y$ and $y_1 = y$, and the equation becomes $x_1^2 + \frac{3}{4}y_1^2 + x_1 - \frac{1}{2}y_1 = 0$.

Now, whenever $\alpha_i \neq 0$, we can replace x_i by $x_i - \beta_i / (2\alpha_i)$, and thereby eliminate the term $\beta_i x_i$ from the equation. This transformation is just a translation, which is also an isometry.

4 Bilinear Maps and Quadratic Forms

For example, suppose we have $x^2 - 3x = 0$. Then we are completing the square again, but this time in one variable. So $x^2 - 3x = 0$ is just $(x - \frac{3}{2})^2 - \frac{9}{4} = 0$ and we use $x_1 = x - \frac{3}{2}$ to write it as $x_1^2 - \frac{9}{4}$.

If $\alpha_i = 0$, then we cannot eliminate the term $\beta_i x_i$. Let us permute the coordinates such that $\alpha_i \neq 0$ for $1 \leq i \leq r$, and $\beta_i \neq 0$ for $r+1 \leq i \leq r+s$.

If $s > 1$, we want to leave the x_i alone for $1 \leq i \leq r$ but replace $\sum_{i=1}^s \beta_{r+i} x_{r+i}$ by βx_{r+1} . To see that we can do this via an orthogonal transformation we use Theorem 4.5.3. Suppose our orthonormal basis was e_1, \dots, e_n . Then we can extend

$$e_1, \dots, e_r, \frac{1}{\sqrt{\sum_{i=1}^s \beta_{r+i}^2}} \sum_{i=1}^s \beta_{r+i} e_{r+i}$$

to an orthonormal basis of our euclidean space. Note that the $r+1$ th vector is chosen so that in this basis our equation will just have the term $(\sqrt{\sum_{i=1}^s \beta_{r+i}^2}) x'_{r+1}$. So we have reduced our equation to at most one non-zero β_i ; either there are no linear terms at all, or there is just β_{r+1} . Given that $(\sqrt{\sum_{i=1}^s \beta_{r+i}^2})$ is positive we might think that our linear term will have a positive coefficient. In fact, by dividing through by a constant we can choose it to be -1 , which we do for convenience.

Finally, if there is a linear term, so $\beta_{r+1} \neq 0$ (and in fact can be thought to be -1 by the above comment), then we can perform the translation that replaces x_{r+1} by $x_{r+1} - \gamma/\beta_{r+1}$, and thereby eliminate the constant γ . When there is no linear term then we divide the equation through by a constant, to assume that γ is 0 or -1 and we put γ on the right hand side for convenience.

We have proved the following theorem:

Theorem 4.8.1. *By rigid motions of euclidean space, we can transform the set defined by the general second degree equation (†) into the set defined by an equation having one of the following three forms:*

$$\begin{aligned} \sum_{i=1}^r \alpha_i x_i^2 &= 0, \\ \sum_{i=1}^r \alpha_i x_i^2 &= 1, \\ \sum_{i=1}^r \alpha_i x_i^2 - x_{r+1} &= 0. \end{aligned}$$

Here $0 \leq r \leq n$ and $\alpha_1, \dots, \alpha_r$ are non-zero constants, and in the third case $r < n$.

We shall assume that $r \neq 0$, because otherwise we have a linear equation. The sets defined by the first two types of equation are called *central quadrics* because they have central symmetry; i.e. if a vector \mathbf{v} satisfies the equation, then so does $-\mathbf{v}$.

We shall now consider the types of curves and surfaces that can arise in the familiar cases $n = 2$ and $n = 3$. These different types correspond to whether the α_i are positive, negative or zero, and whether $\gamma = 0$ or 1.

4 Bilinear Maps and Quadratic Forms

We shall use x, y, z instead of x_1, x_2, x_3 , and α, β, γ instead of $\alpha_1, \alpha_2, \alpha_3$. We shall assume also that α, β, γ are all strictly positive, and write $-\alpha$, etc., for the negative case. When the coefficient of the right hand side is 0, we will divide through by -1 at will. For example, Case (i) in the next section contains both $\alpha x^2 = 0$ and $-\alpha x^2 = 0$, which of course need not be counted twice. Moreover, if swapping the names of x and y (whilst swapping the arbitrary positive real numbers α and β) gives the same equation, we will only consider it once. For example, we do this for the list in the next section by only listing Case (vii) once ($-\alpha x^2 + \beta y^2 = 1$ is also in this case).

4.8.2 The case $n = 2$

When $n = 2$ we have the following possibilities.

- (i) $\alpha x^2 = 0$. This just defines the line $x = 0$ (the y -axis).
- (ii) $\alpha x^2 = 1$. This defines the two parallel lines $x = \pm 1/\sqrt{\alpha}$.
- (iii) $-\alpha x^2 = 1$. This is the empty set!
- (iv) $\alpha x^2 + \beta y^2 = 0$. The single point $(0, 0)$.
- (v) $\alpha x^2 - \beta y^2 = 0$. This defines two straight lines $y = \pm \sqrt{\alpha/\beta} x$, which intersect at $(0, 0)$.
- (vi) $\alpha x^2 + \beta y^2 = 1$. An ellipse.
- (vii) $\alpha x^2 - \beta y^2 = 1$. A hyperbola.
- (viii) $-\alpha x^2 - \beta y^2 = 1$. The empty set again.
- (ix) $\alpha x^2 - y = 0$. A parabola.

4.8.3 The case $n = 3$

When $n = 3$, we still get the nine possibilities (i) – (ix) that we had in the case $n = 2$, but now they must be regarded as equations in the three variables x, y, z that happen not to involve z .

So, in Case (i), we now get the plane $x = 0$, in Case (ii) we get two parallel planes $x = \pm 1/\sqrt{\alpha}$, in Case (iv) we get the line $x = y = 0$ (the z -axis), in Case (v) two intersecting planes $y = \pm \sqrt{\alpha/\beta} x$, and in Cases (vi), (vii) and (ix), we get, respectively, elliptical, hyperbolic and parabolic cylinders.

The remaining cases involve all of x, y and z . We omit $-\alpha x^2 - \beta y^2 - \gamma z^2 = 1$, which is empty.

- (x) $\alpha x^2 + \beta y^2 + \gamma z^2 = 0$. The single point $(0, 0, 0)$.
- (xi) $\alpha x^2 + \beta y^2 - \gamma z^2 = 0$. See Fig. 1.

4 Bilinear Maps and Quadratic Forms

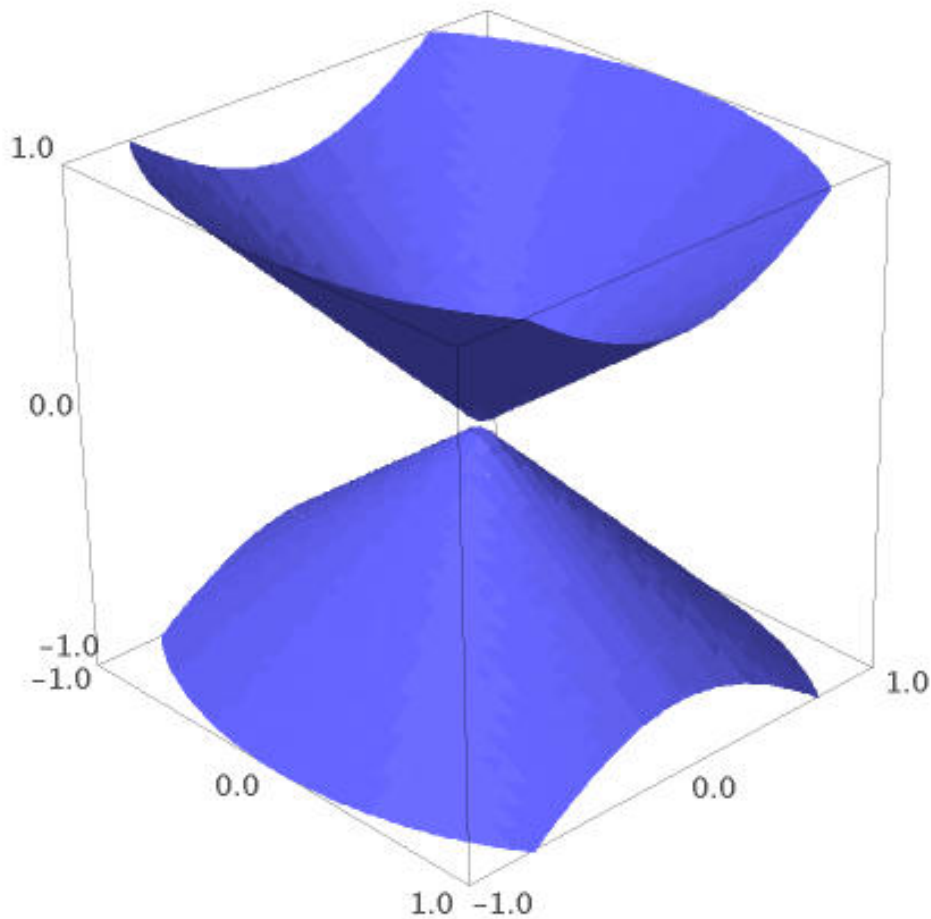


Figure 1: $\frac{1}{2}x^2 + y^2 - z^2 = 0$

This is an elliptical cone. The cross sections parallel to the xy -plane are ellipses of the form $\alpha x^2 + \beta y^2 = c$, whereas the cross sections parallel to the other coordinate planes are generally hyperbolas. Notice also that if a particular point (a, b, c) is on the surface, then so is $t(a, b, c)$ for any $t \in \mathbb{R}$. In other words, the surface contains the straight line through the origin and any of its points. Such lines are called *generators*. When each point of a 3-dimensional surface lies on one or more generators, it is possible to make a model of the surface with straight lengths of wire or string.

(xii) $\alpha x^2 + \beta y^2 + \gamma z^2 = 1$. An ellipsoid. See Fig. 2.

4 Bilinear Maps and Quadratic Forms

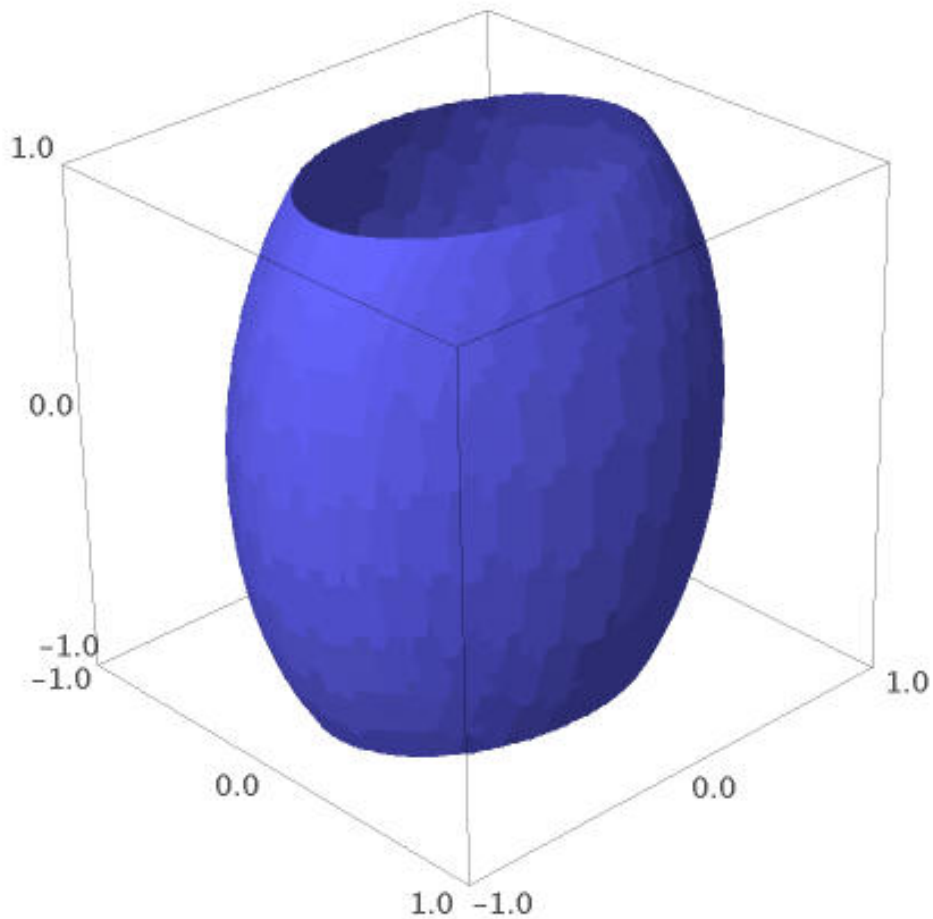


Figure 2: $2x^2 + y^2 + \frac{1}{2}z^2 = 1$

This is a “squashed sphere”. It is bounded, and hence clearly has no generators. Notice that if α , β , and γ are distinct, it has only the finite group of symmetries given by reflections in x , y and z , but if some two of the coefficients coincide, it picks up an infinite group of rotation symmetries.

(xiii) $\alpha x^2 + \beta y^2 - \gamma z^2 = 1$. A hyperboloid. See Fig. 3.

4 Bilinear Maps and Quadratic Forms

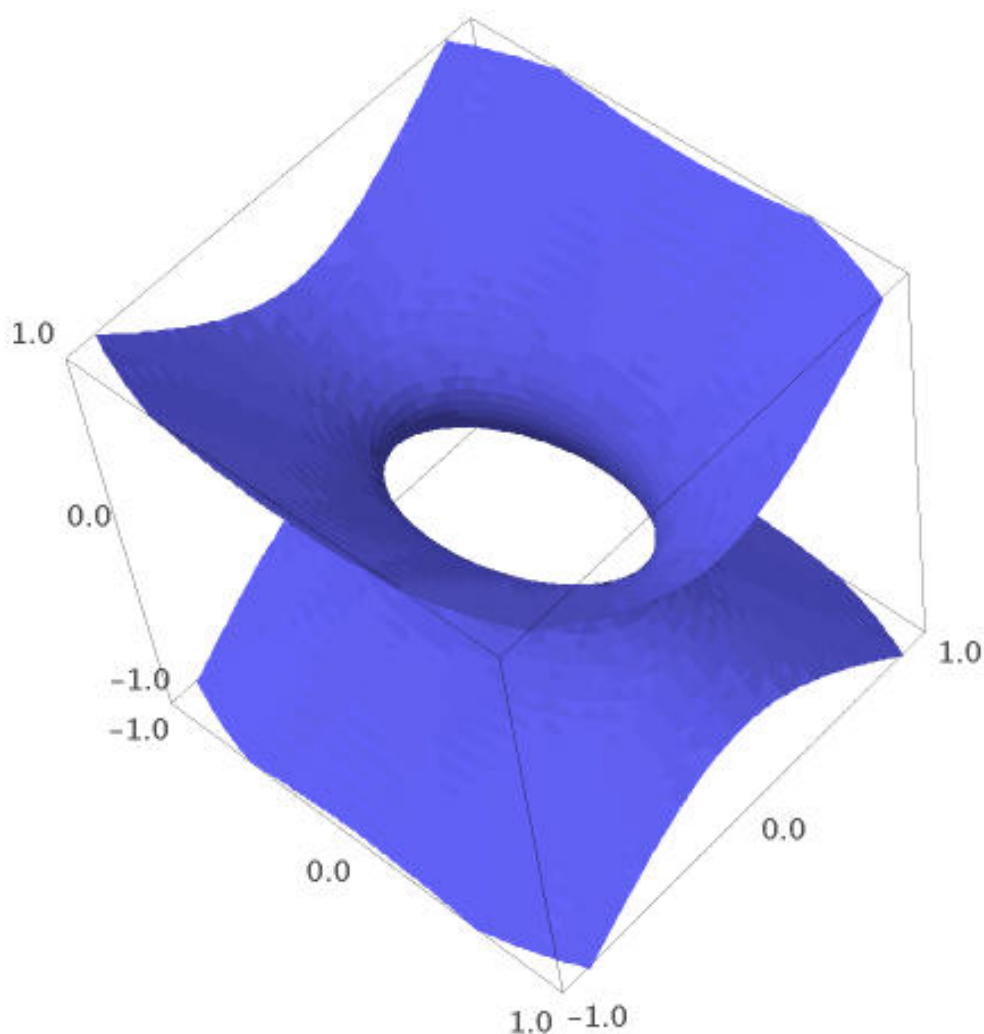


Figure 3: $3x^2 + 8y^2 - 8z^2 = 1$

There are two types of 3-dimensional hyperboloids. This one is connected, and is known as a *hyperboloid of one sheet*. Any cross-section in the xy direction will be an ellipse, and these get larger as z grows (notice the hole in the middle in the picture). Although it is not immediately obvious, each point of this surface lies on exactly two generators; that is, lines that lie entirely on the surface. For each $\lambda \in \mathbb{R}$, the line defined by the pair of equations

$$\sqrt{\alpha}x - \sqrt{\gamma}z = \lambda(1 - \sqrt{\beta}y); \quad \lambda(\sqrt{\alpha}x + \sqrt{\gamma}z) = 1 + \sqrt{\beta}y.$$

lies entirely on the surface; to see this, just multiply the two equations together. The same applies to the lines defined by the pairs of equations

$$\sqrt{\beta}y - \sqrt{\gamma}z = \mu(1 - \sqrt{\alpha}x); \quad \mu(\sqrt{\beta}y + \sqrt{\gamma}z) = 1 + \sqrt{\alpha}x.$$

4 Bilinear Maps and Quadratic Forms

It can be shown that each point on the surface lies on exactly one of the lines in each of these two families.

There is a photo at http://home.cc.umanitoba.ca/~gunderso/model_photos/misc/hyperboloid_of_one_sheet.jpg depicting a rather nice wooden model of a hyperboloid of one sheet, which gives a good idea how these lines sit inside the surface.

(xiv) $ax^2 - \beta y^2 - \gamma z^2 = 1$. Another kind of hyperboloid. See Fig. 4.

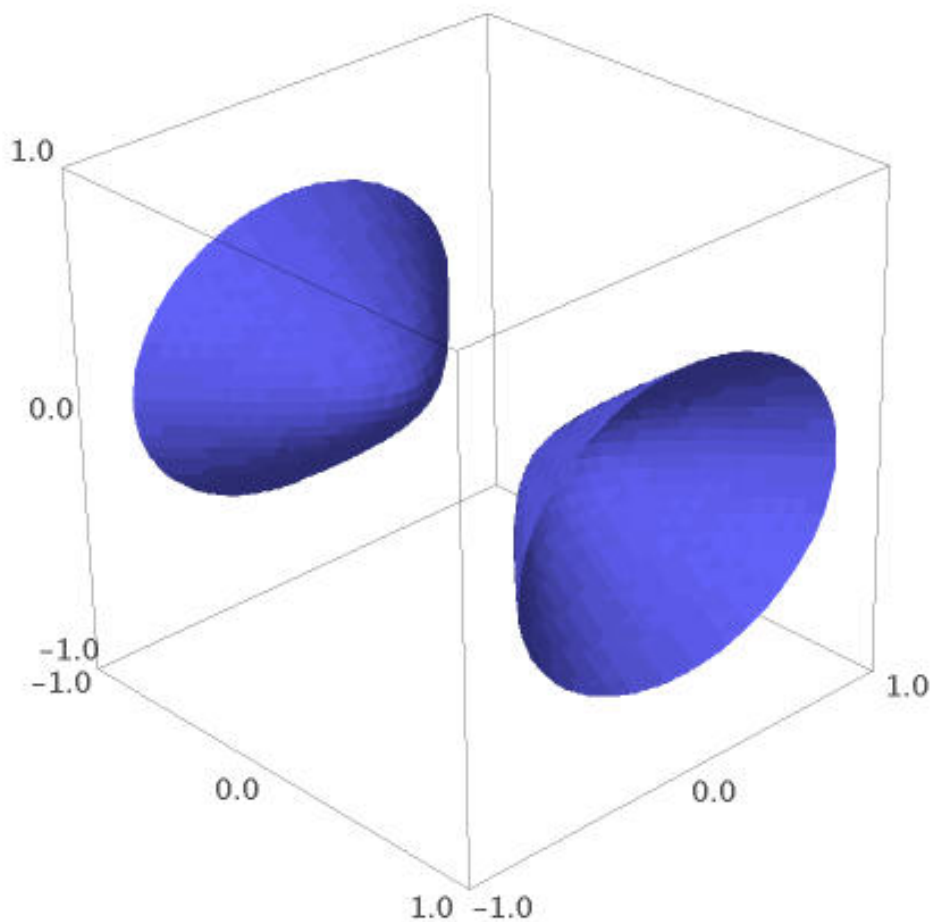


Figure 4: $8x^2 - 12y^2 - 20z^2 = 1$

This one has two connected components and is called a *hyperboloid of two sheets*. It does not have generators.

4 Bilinear Maps and Quadratic Forms

(xv) $\alpha x^2 + \beta y^2 - z = 0$. An elliptical paraboloid. See Fig. 5.

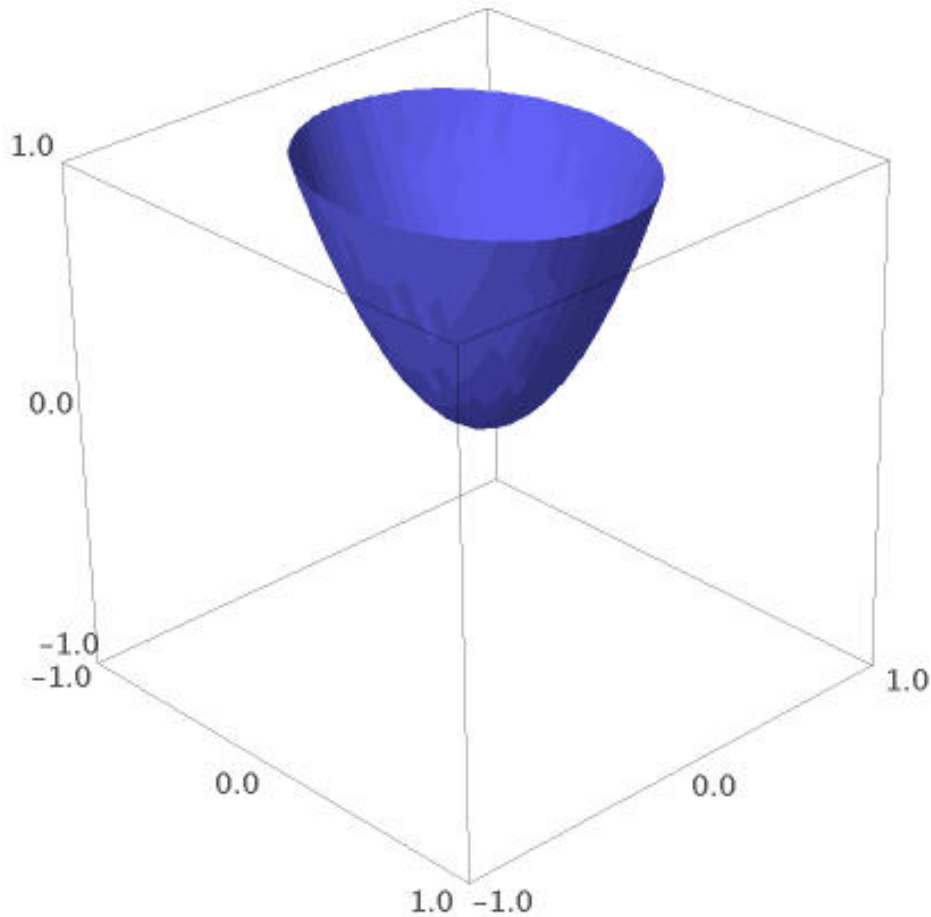


Figure 5: $2x^2 + 3y^2 - z = 0$

Cross-sections of this surface parallel to the xy plane are ellipses, while cross-sections in the yz and xz directions are parabolas. It can be regarded as the limit of a family of hyperboloids of two sheets, where one “cap” remains at the origin and the other recedes to infinity.

(xvi) $\alpha x^2 - \beta y^2 - z = 0$. A hyperbolic paraboloid (a rather elegant saddle shape). See Fig. 6.

4 Bilinear Maps and Quadratic Forms

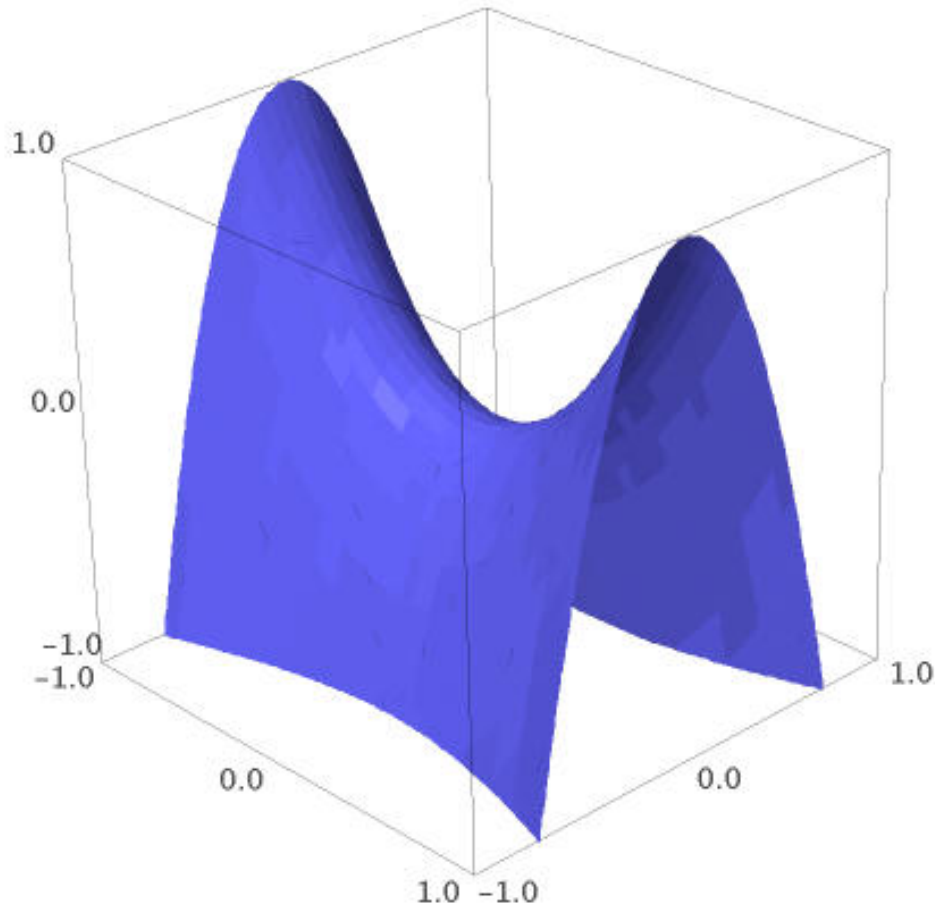


Figure 6: $x^2 - 4y^2 - z = 0$

As in the case of the hyperboloid of one sheet, there are two generators passing through each point of this surface, one from each of the following two families of lines:

$$\begin{aligned} \lambda(\sqrt{\alpha}x - \sqrt{\beta}y) = z; & \quad \sqrt{\alpha}x + \sqrt{\beta}y = \lambda. \\ \mu(\sqrt{\alpha}x + \sqrt{\beta}y) = z; & \quad \sqrt{\alpha}x - \sqrt{\beta}y = \mu. \end{aligned}$$

Just as the elliptical paraboloid was a limiting case of a hyperboloid of two sheets, so the hyperbolic paraboloid is a limiting case of a hyperboloid of one sheet: you can imagine gradually deforming the hyperboloid of one sheet so the elliptical hole in the middle becomes bigger and bigger, and the result is the hyperbolic paraboloid.

4.9 Singular value decomposition

In this section we want to study what linear maps $T : V \rightarrow W$ between euclidean spaces look like? From MA106 Linear Algebra we know that we can choose bases in V and W such that the matrix of T in Smith normal form is $\left(\begin{array}{c|c} I_n & 0 \\ \hline 0 & 0 \end{array} \right)$ where n is the rank of T . This answer is unsatisfactory because it does not take the euclidean geometry of V and W into account. In other words, we want to choose orthonormal bases, not just any bases. This leads us to the singular value decomposition, SVD for short.

Notation: We will see various diagonal matrices in the following so we will use the shorthand $\text{diag}(d_1, \dots, d_n)$ for an $n \times n$ diagonal matrix with diagonal entries d_1, \dots, d_n .

Theorem 4.9.1 (SVD for linear maps). *Suppose $T : V \rightarrow W$ is a linear map of rank n between euclidean spaces. Then there exist unique positive numbers $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_n > 0$, called the singular values of T , and orthonormal bases of V and W such that the matrix of T with respect to these bases is*

$$\left(\begin{array}{c|c} D & 0 \\ \hline 0 & 0 \end{array} \right) \quad \text{where } D = \text{diag}(\gamma_1, \dots, \gamma_n).$$

Proof. We will consider a new bilinear form on V defined as follows.

$$\mathbf{u} \star \mathbf{v} := T(\mathbf{u}) \cdot T(\mathbf{v}).$$

Note that $\mathbf{v} \star \mathbf{v} = T(\mathbf{v}) \cdot T(\mathbf{v}) \geq 0$; we call such a bilinear form *positive semidefinite* (note that it need not be positive definite because T can have a non-zero kernel). By Theorem 4.7.3, there exist unique constants $\alpha_1 \geq \dots \geq \alpha_m$ (eigenvalues of the matrix of the \star bilinear form) and an orthonormal basis $\mathbf{e}_1, \dots, \mathbf{e}_m$ of V such that the bilinear form \star is given by $\text{diag}(\alpha_1, \dots, \alpha_m)$ in this basis. Since \star is positive semidefinite we see that all α_i are non-negative. Suppose $\alpha_k > 0$ is the last positive eigenvalue, that is, $\alpha_{k+1} = \dots = \alpha_m = 0$.

Since $T(\mathbf{e}_i) \cdot T(\mathbf{e}_j) = \mathbf{e}_i \star \mathbf{e}_j = \delta_{ij}\alpha_i$, we deduce that $T(\mathbf{e}_{k+1}) = \dots = T(\mathbf{e}_m) = 0$ (the bilinear form \cdot on W is positive definite) and $T(\mathbf{e}_1), \dots, T(\mathbf{e}_k)$ form an orthogonal set of vectors in W . It follows that k is the rank of T since a set of orthogonal vectors is linearly independent (see the proof of Theorem 4.5.3). Thus, $k = n$. We define $\gamma_i := \sqrt{\alpha_i}$ for all $i \leq k$.

We now use these image vectors $T(\mathbf{e}_i)$ to build an orthonormal basis of W . Since $T(\mathbf{e}_i) \cdot T(\mathbf{e}_i) = \mathbf{e}_i \star \mathbf{e}_i = \alpha_i$, we know that $|T(\mathbf{e}_i)| = \sqrt{\alpha_i} = \gamma_i$. Let $\mathbf{f}_i := \frac{T(\mathbf{e}_i)}{\gamma_i}$ for all $i \leq n$. We can then extend this orthonormal set of vectors to an orthonormal basis of W by the Gram-Schmidt process (Theorem 4.5.3). Since $T(\mathbf{e}_i) = \gamma_i \mathbf{f}_i$ for $i \leq n$ and $T(\mathbf{e}_j) = 0$ for $j > n$, the matrix of T with respect to these bases has the required form.

It remains to prove the uniqueness of the singular values. Suppose we have orthonormal bases $\mathbf{e}'_1, \dots, \mathbf{e}'_m$ of V and $\mathbf{f}'_1, \dots, \mathbf{f}'_s$ of W , in which T is represented by a matrix $\left(\begin{array}{c|c} B & 0 \\ \hline 0 & 0 \end{array} \right)$ where $B = \text{diag}(\beta_1, \dots, \beta_t)$ with $\beta_1 \geq \dots \geq \beta_t > 0$. Put $\beta_i = 0$ for $i > t$. Then $\mathbf{e}'_i \star \mathbf{e}'_j = \beta_i \mathbf{f}'_i \cdot \beta_j \mathbf{f}'_j = \delta_{ij}\beta_i^2$. Thus, $\text{diag}(\beta_1^2, \dots, \beta_m^2)$ is the matrix of the bilinear form \star in the basis $\mathbf{e}'_1, \dots, \mathbf{e}'_m$. Uniqueness in Theorem 4.7.3 implies the uniqueness of the singular values. \square

4 Bilinear Maps and Quadratic Forms

Before we proceed with some examples, all on the standard euclidean spaces \mathbb{R}^n , let us restate the SVD for matrices:

Corollary 4.9.2 (SVD for matrices). *Given any real $k \times m$ matrix A , there exist unique singular values $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_n > 0$ and (non-unique) orthogonal matrices P and Q such that*

$$\left(\begin{array}{c|c} D & 0 \\ \hline 0 & 0 \end{array} \right) = P^T A Q \text{ where } D = \text{diag}(\gamma_1, \dots, \gamma_n).$$

Equivalently, we say the SVD of A is

$$A = P \left(\begin{array}{c|c} D & 0 \\ \hline 0 & 0 \end{array} \right) Q^T \text{ where } D = \text{diag}(\gamma_1, \dots, \gamma_n).$$

Example. Consider a linear map $\mathbb{R}^2 \rightarrow \mathbb{R}^2$, given by the symmetric matrix $A = \begin{pmatrix} 1 & 3 \\ 3 & 1 \end{pmatrix}$, in the example from Section 4.7. There we found the orthogonal matrix $P = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \end{pmatrix}$ with $P^T A P = \begin{pmatrix} 4 & 0 \\ 0 & -2 \end{pmatrix}$. This is not the SVD of A because the diagonal matrix contains a negative entry. To get to the SVD we just need to pick different bases for the domain and the range: the columns $\mathbf{c}_1, \mathbf{c}_2$ can still be a basis of the domain, while the basis of the range could become $\mathbf{c}_1, -\mathbf{c}_2$. This is the SVD:

$$P = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}, Q = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \end{pmatrix}, P^T A Q = \begin{pmatrix} 4 & 0 \\ 0 & 2 \end{pmatrix}.$$

The same method works for any symmetric matrix: the SVD is just orthogonal diagonalisation with additional care needed for signs. If the matrix is not symmetric, we need to follow the proof of Theorem 4.9.1 during the calculation.

Example. Consider a linear map $\mathbb{R}^3 \rightarrow \mathbb{R}^2$, given by $A = \begin{pmatrix} 4 & 11 & 14 \\ 8 & 7 & -2 \end{pmatrix}$. Since $\mathbf{x} \star \mathbf{y} = A\mathbf{x} \cdot A\mathbf{y} = (A\mathbf{x})^T A\mathbf{y} = \mathbf{x}^T (A^T A)\mathbf{y}$, the matrix of the bilinear form \star in the standard basis is

$$A^T A = \begin{pmatrix} 4 & 8 \\ 11 & 7 \\ 14 & -2 \end{pmatrix} \begin{pmatrix} 4 & 11 & 14 \\ 8 & 7 & -2 \end{pmatrix} = \begin{pmatrix} 80 & 100 & 40 \\ 100 & 170 & 140 \\ 40 & 140 & 200 \end{pmatrix}.$$

The eigenvalues of this matrix are 360, 90 and 0. Hence the singular values of A are

$$\gamma_1 = \sqrt{360} = 6\sqrt{10} \geq \gamma_2 = \sqrt{90} = 3\sqrt{10}.$$

At this stage we are assured of the existence of orthogonal matrices P and Q such that

$$P^T A Q = \begin{pmatrix} 6\sqrt{10} & 0 & 0 \\ 0 & 3\sqrt{10} & 0 \end{pmatrix}.$$

4 Bilinear Maps and Quadratic Forms

To find such orthogonal matrices we first need to find an orthonormal basis of eigenvectors of $A^T A$. Since the eigenvalues are distinct on this occasion we only need to find an eigenvector for each eigenvalue and normalise it so it has length 1. This leads to:

$$\mathbf{e}_1 = \begin{pmatrix} 1/3 \\ 2/3 \\ 2/3 \end{pmatrix}, \quad \mathbf{e}_2 = \begin{pmatrix} -2/3 \\ -1/3 \\ 2/3 \end{pmatrix}, \quad \mathbf{e}_3 = \begin{pmatrix} 2/3 \\ -2/3 \\ 1/3 \end{pmatrix}.$$

These make up Q . Then we need to find the images of these vectors under A divided by the corresponding singular value (so only the eigenvectors for the non-zero eigenvalues of $A^T A$):

$$\mathbf{f}_1 = \frac{1}{6\sqrt{10}} A\mathbf{e}_1 = \begin{pmatrix} 3/\sqrt{10} \\ 1/\sqrt{10} \end{pmatrix}, \quad \mathbf{f}_2 = \frac{1}{3\sqrt{10}} A\mathbf{e}_2 = \begin{pmatrix} 1/\sqrt{10} \\ -3/\sqrt{10} \end{pmatrix}.$$

The proof says we need to extend this to a basis of W , which is easy here because we already have two vectors and so we don't need anymore for a basis of \mathbb{R}^2 . Hence, the orthogonal matrices are

$$P = \begin{pmatrix} 3/\sqrt{10} & 1/\sqrt{10} \\ 1/\sqrt{10} & -3/\sqrt{10} \end{pmatrix}, \quad Q = \begin{pmatrix} 1/3 & -2/3 & 2/3 \\ 2/3 & -1/3 & -2/3 \\ 2/3 & 2/3 & 1/3 \end{pmatrix}.$$

4.10 The complex story

The results in Subsection 4.7 applied only to vector spaces over the real numbers \mathbb{R} . There are corresponding results for spaces over the complex numbers \mathbb{C} , which we shall summarize here. We only include one proof, although the others are similar and analogous to those for spaces over \mathbb{R} .

4.10.1 Sesquilinear forms

The key thing that made everything work over \mathbb{R} was the fact that if x_1, \dots, x_n are real numbers, and $x_1^2 + \dots + x_n^2 = 0$, then all the x_i are zero. This doesn't work over \mathbb{C} : take $x_1 = 1$ and $x_2 = i$. But we do have something similar if we bring *complex conjugation* into play. As usual, for $z \in \mathbb{C}$, we let \bar{z} denote the complex conjugate of z . Then if $z_1\bar{z}_1 + \dots + z_n\bar{z}_n = 0$, each z_i must be zero. So we need to "put bars on half of our formulae". Notice that there was a hint of this in the proof of Proposition 4.7.2.

We'll do this as follows.

Definition 4.10.1. A *sesquilinear form* on a complex vector space V is a function $\tau : V \times V \rightarrow \mathbb{C}$ such that

$$\tau(\mathbf{v}, a_1\mathbf{w}_1 + a_2\mathbf{w}_2) = a_1\tau(\mathbf{v}, \mathbf{w}_1) + a_2\tau(\mathbf{v}, \mathbf{w}_2)$$

(as before), but

$$\tau(a_1\mathbf{v}_1 + a_2\mathbf{v}_2, \mathbf{w}) = \bar{a}_1\tau(\mathbf{v}_1, \mathbf{w}) + \bar{a}_2\tau(\mathbf{v}_2, \mathbf{w}),$$

for all vectors v_1, v_2, v, w_1, w_2, w and all $a_1, a_2 \in \mathbb{C}$.

4 Bilinear Maps and Quadratic Forms

We say such a form is *hermitian symmetric* if

$$\tau(\mathbf{w}, \mathbf{v}) = \overline{\tau(\mathbf{v}, \mathbf{w})}.$$

The word “sesquilinear” literally means “one-and-a-half-times-linear” from its Latin meaning – it’s linear in the second argument, but only halfway there in the first argument! We’ll often abbreviate “hermitian-symmetric sesquilinear form” to just “hermitian form”.

We can represent these by matrices in a similar way to bilinear forms. If τ is a sesquilinear form, and $\mathbf{e}_1, \dots, \mathbf{e}_n$ is a basis of V , we define the matrix of τ to be the matrix A whose i, j entry is $\tau(\mathbf{e}_i, \mathbf{e}_j)$. Then we have

$$\tau(\mathbf{v}, \mathbf{w}) = (\underline{\mathbf{v}}^T) A \underline{\mathbf{w}}$$

where $\underline{\mathbf{v}}$ and $\underline{\mathbf{w}}$ are the coordinates of \mathbf{v} and \mathbf{w} as usual. We’ll shorten this to $\underline{\mathbf{v}}^* A \underline{\mathbf{w}}$, where the $*$ denotes “conjugate transpose”. The condition to be hermitian symmetric translates to the relation $a_{ji} = \overline{a_{ij}}$, so τ is hermitian if and only if A satisfies $A^* = A$.

We have a version here of Sylvester’s two theorems (Proposition 4.4.3 and Theorem 4.4.4):

Theorem 4.10.2. *If τ is a hermitian form on a complex vector space V , there is a basis of V in which the matrix of τ is given by*

$$\left(\begin{array}{c|c|c} I_t & & \\ \hline & -I_u & \\ \hline & & 0 \end{array} \right)$$

for some uniquely determined integers t and u .

As in the real case, we call the pair (t, u) the *signature* of τ , and we say τ is *positive definite* if its signature is $(n, 0)$ (if V is an n -dimensional space). In this case, the theorem tells us that there is a basis of V in which the matrix of τ is the identity, and in such a basis we have

$$\tau(\mathbf{v}, \mathbf{v}) = \sum_{i=1}^n |v_i|^2$$

where v_1, \dots, v_n are the coordinates of \mathbf{v} . Hence $\tau(\mathbf{v}, \mathbf{v}) > 0$ for all non-zero $\mathbf{v} \in V$.

Just as we defined a euclidean space to be a real vector space with a choice of positive definite bilinear form, we have a similar definition here:

Definition 4.10.3. *A Hilbert space is a finite-dimensional complex vector space endowed with a choice of positive-definite hermitian-symmetric sesquilinear form.*

These are the complex analogues of euclidean spaces. If V is a Hilbert space, we write $\mathbf{v} \cdot \mathbf{w}$ for the sesquilinear form on V , and we refer to it as an *inner product*. For any Hilbert space, we can always find a basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ of V such that $\mathbf{e}_i \cdot \mathbf{e}_j = \delta_{ij}$ (an orthonormal basis). Then we can write the inner product matrix-wise as

$$\mathbf{v} \cdot \mathbf{w} = \underline{\mathbf{v}}^* \underline{\mathbf{w}},$$

4 Bilinear Maps and Quadratic Forms

where $\underline{\mathbf{v}}$ and $\underline{\mathbf{w}}$ are the coordinates of \mathbf{v} and \mathbf{w} and $\underline{\mathbf{v}}^* = \overline{\underline{\mathbf{v}}^T}$ as before.

The canonical example of a Hilbert space is \mathbb{C}^n , with the standard inner product given by

$$\mathbf{v} \cdot \mathbf{w} = \sum_{i=1}^n \overline{v_i} w_i,$$

for which the standard basis is obviously orthonormal.

Remark. *Technically, we should say “finite-dimensional Hilbert space”. There are lots of interesting infinite-dimensional Hilbert spaces, but we won’t say anything about them in this course. (Curiously, one never seems to come across infinite-dimensional euclidean spaces.)*

4.10.2 Operators on Hilbert spaces

In our study of linear operators on euclidean spaces, the idea of the *adjoint* of an operator was important. There’s an analogue of it here:

Definition 4.10.4. Let $T : V \rightarrow V$ be a linear operator on a Hilbert space V . Then there is a unique linear operator $T^* : V \rightarrow V$ (the *hermitian adjoint* of T) such that

$$T(\mathbf{v}) \cdot \mathbf{w} = \mathbf{v} \cdot T^*(\mathbf{w}).$$

It’s clear that if A is the matrix of T in an orthonormal basis, then the matrix of T^* is A^* .

Definition 4.10.5. We say that T is

- *selfadjoint* if $T^* = T$,
- *unitary* if $T^* = T^{-1}$,
- *normal* if $T^*T = TT^*$.

Exercise. If T is unitary, then $T(\mathbf{u}) \cdot T(\mathbf{v}) = \mathbf{u} \cdot \mathbf{v}$ for all \mathbf{u}, \mathbf{v} in V .

Using this exercise we can also replicate Proposition 4.6.5 in the complex world. This shows that ‘unitary’ is the complex analogue of ‘orthogonal’. The proof is entirely similar to that of Proposition 4.6.5 (which comes before the statement).

Proposition 4.10.6. *Let e_1, \dots, e_n be an orthonormal basis of a Hilbert space V . A linear map T is unitary if and only if $T(e_1), \dots, T(e_n)$ is an orthonormal basis of V .*

If A is the matrix of T in an orthonormal basis, then it’s clear that T is selfadjoint if and only if $A^* = A$ (a hermitian-symmetric matrix), unitary if and only if $A^* = A^{-1}$ (a *unitary matrix*), and normal if and only if $A^*A = AA^*$ (a *normal matrix*). In other words, these properties are preserved under unitary base changes:

Lemma 4.10.7. *If $A \in \mathbb{C}^{n,n}$ is normal (selfadjoint, unitary) and $P \in \mathbb{C}^{n,n}$ is unitary, then P^*AP is normal (selfadjoint, unitary).*

4 Bilinear Maps and Quadratic Forms

Proof. Let $B = P^*AP$. Using the property $(MN)^* = N^*M^*$, we compute that in the first (normal) case,

$$BB^* = (P^*AP)(P^*AP)^* = P^*APP^*A^*P = P^*AA^*P = P^*A^*AP = (P^*A^*P)(P^*AP) = B^*B.$$

In the second (selfadjoint) case, $B^* = P^*A^*P = P^*AP = B$. In the third (unitary) case, $BB^* = P^*APP^*A^*P = P^*AA^*P = P^*P = I$. \square

Notice that if A is unitary and the entries of A are real, then A must be orthogonal, but the definition also includes things like

$$\begin{pmatrix} i & 0 \\ 0 & i \end{pmatrix}.$$

Similarly, a matrix with real entries is hermitian-symmetric if and only if it's symmetric, but

$$\begin{pmatrix} 2 & i \\ -i & 3 \end{pmatrix}$$

is a hermitian-symmetric matrix that's not symmetric.

Both selfadjoint and unitary operators are normal. The generalisation of Theorem 4.7.3 applies to all three types of operators.

Theorem 4.10.8. *The following statements hold for a linear operator $T : V \rightarrow V$ on a Hilbert space.*

- (i) *T is normal if and only if there exists an orthonormal basis of V consisting of eigenvectors of T .*
- (ii) *T is selfadjoint if and only if there exists an orthonormal basis of V consisting of eigenvectors of T with real eigenvalues.*
- (iii) *T is unitary if and only if there exists an orthonormal basis of V consisting of eigenvectors of T with eigenvalues of absolute value 1.*

Example. Let $A = \begin{pmatrix} 6 & 2+2i \\ 2-2i & 4 \end{pmatrix}$. Then

$$c_A(x) = (6-x)(4-x) - (2+2i)(2-2i) = x^2 - 10x + 16 = (x-2)(x-8),$$

so the eigenvalues are 2 and 8. Corresponding eigenvectors are $\mathbf{v}_1 = (1+i, -2)^T$ and $\mathbf{v}_2 = (1+i, 1)^T$. We find that $|\mathbf{v}_1|^2 = \mathbf{v}_1^*\mathbf{v}_1 = 6$ and $|\mathbf{v}_2|^2 = 3$, so we divide by their lengths to get an orthonormal basis $\mathbf{v}_1/|\mathbf{v}_1|, \mathbf{v}_2/|\mathbf{v}_2|$ of \mathbb{C}^2 . Then the matrix

$$P = \begin{pmatrix} \frac{1+i}{\sqrt{6}} & \frac{1+i}{\sqrt{3}} \\ \frac{-2}{\sqrt{6}} & \frac{1}{\sqrt{3}} \end{pmatrix}$$

having this basis as columns is selfadjoint and satisfies $P^*AP = \begin{pmatrix} 2 & 0 \\ 0 & 8 \end{pmatrix}$.

5 Finitely Generated Abelian Groups

In the first four sections of the course, we've always been thinking about vector spaces over *fields*. The idea of this section is to show that some of the same ideas work with the field K replaced by the integers \mathbb{Z} , even though \mathbb{Z} isn't a field; and that this is strongly related to the *group theory* which most of you will have seen in MA136 Introduction to Abstract Algebra last year. Do not worry if you did not take that module, we will cover all of the group theory we need in the following sections.

5.1 Definitions

Definition 5.1.1. An abelian group is a set G together with a binary operation, which we write as addition, and which satisfies the following properties:

- (i) (*Closure*) for all $g, h \in G$, $g + h \in G$;
- (ii) (*Associativity*) for all $g, h, k \in G$, $(g + h) + k = g + (h + k)$;
- (iii) there exists an element $0_G \in G$ such that:
 - (a) (*Identity*) for all $g \in G$, $g + 0_G = g$; and
 - (b) (*Inverse*) for all $g \in G$ there exists $-g \in G$ such that $g + (-g) = 0_G$;
- (iv) (*Commutativity*) for all $g, h \in G$, $g + h = h + g$.

Usually we just write 0 rather than 0_G . We only write 0_G if we need to distinguish between the zero elements of different groups.

The commutativity axiom (iv) is not part of the definition of a general group, and for general (non-abelian) groups, it is more usual to use multiplicative rather than additive notation. All groups in this course should be assumed to be abelian, although many of the definitions in this section apply equally well to general groups.

Examples. 1. The integers \mathbb{Z} .

2. Fix a positive integer $n > 0$ and let

$$\mathbb{Z}_n = \{0, 1, 2, \dots, n-1\} = \{x \in \mathbb{Z} \mid 0 \leq x < n\}.$$

where addition is computed modulo n . So, for example, when $n = 9$, we have $2 + 5 = 7$, $3 + 8 = 2$, $6 + 7 = 4$, etc. Note that the inverse $-x$ of $x \in \mathbb{Z}_n$ is equal to $n - x$ in this example.

3. Examples from linear algebra. Let K be a field.
- (i) The elements of K form an abelian group under addition.
 - (ii) The non-zero elements of K form an abelian group K^\times under multiplication.
 - (iii) The vectors in any vector space form an abelian group under addition.

Proposition 5.1.2 (The cancellation law). *Let G be any group, and let $g, h, k \in G$. Then $g + h = g + k \Rightarrow h = k$.*

5 Finitely Generated Abelian Groups

Proof. Add $-g$ to both sides of the equation and use the Associativity and Identity axioms. \square

For any group G , $g \in G$, and integer $n > 0$, we define ng to be $g + g + \cdots + g$, with n occurrences of g in the sum. So, for example, $1g = g$, $2g = g + g$, $3g = g + g + g$, etc. We extend this notation to all $n \in \mathbb{Z}$ by defining $0g = 0$ and $(-n)g = -(ng)$ for $-n < 0$. Overall, this defines a scalar action $\mathbb{Z} \times G \rightarrow G$ which allows us to think of abelian groups as “vector spaces over \mathbb{Z} ” (or using precise terminology \mathbb{Z} -modules - algebraic modules will play a significant role in *Rings and Modules* in year 3).

Definition 5.1.3. A group G is called *cyclic* if there exists an element $x \in G$ such that every element of G is of the form mx for some $m \in \mathbb{Z}$.

The element x in the definition is called a *generator* of G . Note that \mathbb{Z} and \mathbb{Z}_n are cyclic with generator $x = 1$.

Definition 5.1.4. A bijection $\phi : G \rightarrow H$ between two (abelian) groups is called an *isomorphism* if $\phi(g + h) = \phi(g) + \phi(h)$ for all $g, h \in G$, and the groups G and H are called *isomorphic* if there is an isomorphism between them.

The notation $G \cong H$ means that G is isomorphic to H ; isomorphic groups are often thought of as being essentially the same group, but with elements having different names.

Note (exercise) that any isomorphism must satisfy $\phi(0_G) = 0_H$ and $\phi(-g) = -\phi(g)$ for all $g \in G$.

Proposition 5.1.5. Any cyclic group G is isomorphic either to \mathbb{Z} or to \mathbb{Z}_n for some $n > 0$.

Proof. Let G be cyclic with generator x . So $G = \{mx \mid m \in \mathbb{Z}\}$. Suppose first that the elements mx for $m \in \mathbb{Z}$ are all distinct. Then the map $\phi : \mathbb{Z} \rightarrow G$ defined by $\phi(m) = mx$ is a bijection, and it is straightforward to check that it is an isomorphism.

Otherwise, we have $lx = mx$ for some $l < m$, and so $(m - l)x = 0$ with $m - l > 0$. Let n be the least integer with $n > 0$ and $nx = 0$. Then the elements $0x = 0, 1x, 2x, \dots, (n - 1)x$ of G are all distinct, because otherwise we could find a smaller n . Furthermore, for any $mx \in G$, we can write $m = rn + s$ for some $r, s \in \mathbb{Z}$ with $0 \leq s < n$. Then $mx = (rn + s)x = sx$, so $G = \{0, 1x, 2x, \dots, (n - 1)x\}$, and the map $\phi : \mathbb{Z}_n \rightarrow G$ defined by $\phi(m) = mx$ for $0 \leq m < n$ is a bijection, and we check that it is an isomorphism. \square

Definition 5.1.6. For an element $g \in G$, the least integer $n > 0$ with $ng = 0$, if it exists, is called the *order* of g and we denote the order of g by $|g|$. If there is no such n , then g has infinite order and we write $|g| = \infty$.

Exercise. If $\phi : G \rightarrow H$ is an isomorphism then $|g| = |\phi(g)|$ for all $g \in G$.

Definition 5.1.7. A group G is *generated* or *spanned* by a subset X of G if every $g \in G$ can be written as a finite sum $\sum_{i=1}^k m_i x_i$, with $m_i \in \mathbb{Z}$ and $x_i \in X$. It is *finitely generated* if it has a finite generating set $X = \{x_1, \dots, x_n\}$.

So a group is cyclic if and only if it has a generating set X with $|X| = 1$.

5 Finitely Generated Abelian Groups

In general, if G is generated by X , then we write $G = \langle X \rangle$ or $G = \langle x_1, \dots, x_n \rangle$ when $X = \{x_1, \dots, x_n\}$ is finite.

Definition 5.1.8. The direct sum of groups G_1, \dots, G_n is defined to be the set $\{(g_1, g_2, \dots, g_n) \mid g_i \in G_i\}$ with component-wise addition

$$(g_1, g_2, \dots, g_n) + (h_1, h_2, \dots, h_n) = (g_1 + h_1, g_2 + h_2, \dots, g_n + h_n).$$

This is a group with identity element $(0, 0, \dots, 0)$ and $-(g_1, g_2, \dots, g_n) = (-g_1, -g_2, \dots, -g_n)$.

In general (non-abelian) group theory this is more often known as the direct product of groups.

The main result of this section, known as the *fundamental theorem of finitely generated abelian groups*, is that every finitely generated abelian group is isomorphic to a direct sum of cyclic groups. (This is not true in general for abelian groups, such as the additive group \mathbb{Q} of rational numbers, which are not finitely generated.)

5.2 Subgroups, cosets and quotient groups

Definition 5.2.1. A subset H of a group G is called a *subgroup* of G if it forms a group under the same operation as that of G .

Lemma 5.2.2. If H is a subgroup of G , then the identity element 0_H of H is equal to the identity element 0_G of G .

Proof. Using the identity axioms for H and G , $0_H + 0_H = 0_H = 0_H + 0_G$. Now by the cancellation law, $0_H = 0_G$. □

The definition of a subgroup is *semantic* in its nature. While it precisely pinpoints what a subgroup is, it is quite cumbersome to use. The following proposition gives a usable criterion.

Proposition 5.2.3. Let H be a subset of a group G . The following statements are equivalent.

- (i) H is a subgroup of G .
- (ii) (a) H is nonempty; and
 - (b) $h_1, h_2 \in H \Rightarrow h_1 + h_2 \in H$; and
 - (c) $h \in H \Rightarrow -h \in H$.
- (iii) (a) H is nonempty; and
 - (b) $h_1, h_2 \in H \Rightarrow h_1 - h_2 \in H$.

Proof. If H is a subgroup of G then it is nonempty as it contains 0 . Moreover, $h_1 - h_2 = h_1 + (-h_2) \in H$ if h_1 and h_2 are in H . Thus, (i) implies (iii).

To show that (iii) implies (ii) we pick $x \in H$. Then $0 = x - x \in H$. Now $-h = 0 - h \in H$ for any $h \in H$. Finally, $h_1 + h_2 = h_1 - (-h_2) \in H$ for all $h_1, h_2 \in H$.

5 Finitely Generated Abelian Groups

To show that (ii) implies (i) we need to verify the four group axioms in H . Two of these, 'Closure', and 'Inverse', are the conditions (b) and (c). The other two axioms are 'Associativity' and 'Identity'. Associativity holds because it holds in G , and H is a subset of G . Since we are assuming that H is nonempty, there exists $h \in H$, and then $-h \in H$ by (c), and $h + (-h) = 0 \in H$ by (b), and so 'Identity' holds, and H is a subgroup. \square

Examples. 1. There are two standard subgroups of any group G : the whole group G itself, and the *trivial* subgroup $\{0\}$ consisting of the identity alone. Subgroups other than G are called *proper* subgroups, and subgroups other than $\{0\}$ are called *non-trivial* subgroups.

2. If g is any element of any group G , then the set of all integer multiples $\{mg \mid m \in \mathbb{Z}\}$ forms a subgroup of G called the cyclic subgroup generated by g .

Let us look at a few specific examples. If $G = \mathbb{Z}$, then $5\mathbb{Z}$, which consists of all multiples of 5, is the cyclic subgroup generated by 5. Of course, we can replace 5 by any integer here, but note that the cyclic groups generated by 5 and -5 are the same.

If $G = \langle g \rangle$ is a finite cyclic group of order n and m is a positive integer dividing n , then the cyclic subgroup generated by mg has order n/m and consists of the elements kmg for $0 \leq k < n/m$.

Exercise. What is the order of the cyclic subgroup generated by mg for general m (where we drop the assumption that $m|n$)?

Exercise. Show that the group of non-zero complex numbers \mathbb{C}^\times under the operation of multiplication has finite cyclic subgroups of all possible orders.

Definition 5.2.4. Let $g \in G$. Then the *coset* $H + g$ is the subset $\{h + g \mid h \in H\}$ of G .

(*Note:* Since our groups are abelian, we have $H + g = g + H$, but in general group theory the right and left cosets $H + g$ and $g + H$ can be different.)

Examples. 1. $G = \mathbb{Z}$, $H = 5\mathbb{Z}$. There are just 5 distinct cosets $H = H + 0 = \{5n \mid n \in \mathbb{Z}\}$, $H + 1 = \{5n + 1 \mid n \in \mathbb{Z}\}$, $H + 2$, $H + 3$, $H + 4$. Note that $H + i = H + j$ whenever $i \equiv j \pmod{5}$.

2. $G = \mathbb{Z}_6$, $H = \{0, 3\}$. There are 3 distinct cosets, $H = H + 3 = \{0, 3\}$, $H + 1 = H + 4 = \{1, 4\}$, and $H + 2 = H + 5 = \{2, 5\}$,

3. $G = \mathbb{C}^\times$, the group of non-zero complex numbers under multiplication and the subgroup $S^1 = \{z, |z| = 1\}$, which is the unit circle. The cosets are circles. There are uncountably many distinct cosets, one for each positive real number (radius of a circle).

Proposition 5.2.5. *The following are equivalent for $g, k \in G$:*

(i) $k \in H + g$.

(ii) $H + g = H + k$.

(iii) $k - g \in H$.

Proof. Clearly $H + g = H + k \Rightarrow k \in H + g$, so (ii) \Rightarrow (i).

5 Finitely Generated Abelian Groups

If $k \in H + g$, then $k = h + g$ for some fixed $h \in H$, so $g = k - h$. Let $f \in H + g$. Then, for some $h_1 \in H$, we have $f = h_1 + g = h_1 + k - h \in H + k$, so $H + g \subseteq H + k$. Similarly, if $f \in H + k$, then for some $h_1 \in H$, we have $f = h_1 + k = h_1 + h + g \in H + g$, so $H + k \subseteq H + g$. Thus $H + g = H + k$, and we have proved that (i) \Rightarrow (ii).

If $k \in H + g$, then, as above, $k = h + g$, so $k - g = h \in H$ and (i) \Rightarrow (iii).

Finally, if $k - g \in H$, then putting $h = k - g$, we have $h + g = k$, so $k \in H + g$, proving (iii) \Rightarrow (i). \square

Corollary 5.2.6. *Two cosets $H + g_1$ and $H + g_2$ of H in G are either equal or disjoint.*

Proof. If $H + g_1$ and $H + g_2$ are not disjoint, then there exists an element $k \in (H + g_1) \cap (H + g_2)$, but then $H + g_1 = H + k = H + g_2$ by Proposition 5.2.5. \square

Corollary 5.2.7. *The cosets of H in G partition G .*

Proposition 5.2.8. *If H is finite, then all cosets have exactly $|H|$ elements.*

Proof. Since $h_1 + g = h_2 + g \Rightarrow h_1 = h_2$ by the cancellation law, it follows that the map $\phi : H \rightarrow H + g$ defined by $\phi(h) = h + g$ is a bijection, and the result follows. \square

Corollary 5.2.7 and Proposition 5.2.8 together imply:

Theorem 5.2.9 (Lagrange's Theorem). *Let G be a finite (abelian) group and H a subgroup of G . Then the order of H divides the order of G .*

Definition 5.2.10. The number of distinct right cosets of H in G is called the *index* of H in G and is written as $|G : H|$.

If G is finite, then we clearly have $|G : H| = |G|/|H|$. But, from the example $G = \mathbb{Z}$, $H = 5\mathbb{Z}$ above, we see that $|G : H|$ can be finite even when G and H are infinite.

Proposition 5.2.11. *Let G be a finite (abelian) group. Then for any $g \in G$, the order $|g|$ of g divides the order $|G|$ of G .*

Proof. Let $|g| = n$. We saw in Example 2 above that the integer multiples $\{mg \mid m \in \mathbb{Z}\}$ of g form a subgroup H of G . By minimality of n , the distinct elements of H are $\{0, g, 2g, \dots, (n-1)g\}$, so $|H| = n$ and the result follows from Lagrange's Theorem. \square

As an application, we can now immediately classify all finite (abelian) groups whose order is prime.

Proposition 5.2.12. *Let G be a (abelian) group having prime order p . Then G is cyclic; that is, $G \cong \mathbb{Z}_p$.*

Proof. Let $g \in G$ with $0 \neq g$. Then $|g| > 1$, but $|g|$ divides p by Proposition 5.2.11, so $|g| = p$. But then G must consist entirely of the integer multiples mg ($0 \leq m < p$) of g , so G is cyclic. \square

5 Finitely Generated Abelian Groups

Definition 5.2.13. If A and B are subsets of a group G , then we define their sum $A + B = \{a + b \mid a \in A, b \in B\}$.

Lemma 5.2.14. If H is a subgroup of the abelian group G and $H + g, H + k$ are cosets of H in G , then $(H + g) + (H + k) = H + (g + k)$.

Proof. Since G is abelian, this follows directly from commutativity and associativity. □

Theorem 5.2.15. Let H be a subgroup of an abelian group G . Then the set G/H of cosets $H + g$ of H in G forms a group under addition of subsets.

Proof. We have just seen that $(H + g) + (H + k) = H + (g + k)$, so we have closure, and associativity follows easily from associativity of G . Since $(H + 0) + (H + g) = H + g$ for all $g \in G$, $H = H + 0$ is an identity element, and since $(H - g) + (H + g) = H - g + g = H$, $H - g$ is an inverse to $H + g$ for all cosets $H + g$. Thus the four group axioms are satisfied and G/H is a group. □

Definition 5.2.16. The group G/H is called the *quotient group* (or the *factor group*) of G by H .

Notice that if G is finite, then $|G/H| = |G : H| = |G|/|H|$. So, although the quotient group seems a rather complicated object at first sight, it is actually a smaller group than G .

Examples. 1. Let $G = \mathbb{Z}$ and $H = m\mathbb{Z}$ for some $m > 0$. Then there are exactly m distinct cosets, $H, H + 1, \dots, H + (m - 1)$. If we add together k copies of $H + 1$, then we get $H + k$. So G/H is cyclic of order m and with generator $H + 1$. So by Proposition 5.1.5, $\mathbb{Z}/m\mathbb{Z} \cong \mathbb{Z}_m$.

2. $G = \mathbb{R}$ and $H = \mathbb{Z}$. The quotient group G/H is isomorphic to the circle subgroup S^1 of the multiplicative group \mathbb{C}^\times . One writes an explicit isomorphism $\phi : G/H \rightarrow S^1$ by $\phi(x + \mathbb{Z}) = e^{2\pi xi}$.

5.3 Homomorphisms and the first isomorphism theorem

Definition 5.3.1. Let G and H be groups. A *homomorphism* ϕ from G to H is a map $\phi : G \rightarrow H$ such that $\phi(g_1 + g_2) = \phi(g_1) + \phi(g_2)$ for all $g_1, g_2 \in G$.

Homomorphisms correspond to linear transformations between vector spaces.

Note that an isomorphism is just a bijective homomorphism. There are two other types of ‘morphism’ that are worth mentioning at this stage.

A homomorphism ϕ is injective if it is an injection; that is, if $\phi(g_1) = \phi(g_2) \Rightarrow g_1 = g_2$. A homomorphism ϕ is surjective if it is a surjection; that is, if $\text{im}(\phi) = H$. Sometimes, a surjective homomorphism is called *epimorphism* while an injective homomorphism is called *monomorphism* but we will not use this terminology in these lectures.

Lemma 5.3.2. Let $\phi : G \rightarrow H$ be a homomorphism. Then $\phi(0_G) = 0_H$ and $\phi(-g) = -\phi(g)$ for all $g \in G$.

5 Finitely Generated Abelian Groups

Proof. Exercise. (Similar to results for linear transformations.) □

Example. Let G be any group, and let $n \in \mathbb{Z}$. Then $\phi : G \rightarrow G$ defined by $\phi(g) = ng$ for all $g \in G$ is a homomorphism.

Kernels and images are defined as for linear transformations of vector spaces.

Definition 5.3.3. Let $\phi : G \rightarrow H$ be a homomorphism. Then the *kernel* $\ker(\phi)$ of ϕ is defined to be the set of elements of G that map onto 0_H ; that is,

$$\ker(\phi) = \{g \mid g \in G, \phi(g) = 0_H\}.$$

Note that by Lemma 5.3.2 above, $\ker(\phi)$ always contains 0_G .

Proposition 5.3.4. Let $\phi : G \rightarrow H$ be a homomorphism. Then ϕ is injective if and only if $\ker(\phi) = \{0_G\}$.

Proof. Since $0_G \in \ker(\phi)$, if ϕ is injective then we must have $\ker(\phi) = \{0_G\}$. Conversely, suppose that $\ker(\phi) = \{0_G\}$, and let $g_1, g_2 \in G$ with $\phi(g_1) = \phi(g_2)$. Then $0_H = \phi(g_1) - \phi(g_2) = \phi(g_1 - g_2)$ (by Lemma 5.3.2), so $g_1 - g_2 \in \ker(\phi)$ and hence $g_1 - g_2 = 0_G$ and $g_1 = g_2$. So ϕ is injective. □

Theorem 5.3.5. Let $\phi : G \rightarrow H$ be a homomorphism. Then $\ker(\phi)$ is a subgroup of G and $\text{im}(\phi)$ is a subgroup of H . Furthermore, if K is a subgroup of a group G then the map $\phi : G \rightarrow G/K$ defined by $\phi(g) = K + g$ is a surjective homomorphism with kernel K .

Proof. The first statement is straightforward using Proposition 5.2.3. For the second, it is clear that ϕ is surjective, and $\phi(g) = 0_{G/K} \Leftrightarrow K + g = K + 0_G \Leftrightarrow g \in K$, so $\ker(\phi) = K$. □

The following lemma explains a connection between quotients and homomorphisms. It clarifies the trickiest point in the proof of the forthcoming *First Isomorphism Theorem*.

Lemma 5.3.6. Let $\phi : G \rightarrow H$ be a group homomorphism with kernel K and let A be a subgroup of G . The induced map $\bar{\phi} : G/A \rightarrow H$ via $\bar{\phi}(A + g) = \phi(g)$ for all $g \in G$ is a group homomorphism if and only if $A \leq K$.

Proof. We need to show that $\bar{\phi}(A + g) = \phi(g)$ does actually define a map $\bar{\phi} : G/A \rightarrow H$. This is not immediately obvious because cosets have different representatives, i.e. we can have $A + g = A + h$ with $g \neq h$. So for $\bar{\phi}$ to make sense we need to ensure $\phi(g) = \phi(h)$ whenever $A + g = A + h$. This is called checking that $\bar{\phi}$ is *well-defined*. Now we know what to check, it is not too difficult. Suppose that $A + g = A + h$. Then $g = a + h$ for some $a \in A$. Hence, $\bar{\phi}$ is well-defined if and only if $\phi(g) = \phi(a) + \phi(h) = \phi(h)$ for all $g, h \in G, a \in A$ if and only if $\phi(a) = 0$ for all $a \in A$ if and only if $A \leq K$.

Now we know that $\bar{\phi}$ is well-defined, we can show it is a homomorphism since ϕ is:

$$\bar{\phi}(A + h) + \bar{\phi}(A + g) = \phi(h) + \phi(g) = \phi(h + g) = \bar{\phi}(A + h + g) = \bar{\phi}((A + h) + (A + g)).$$

□

5 Finitely Generated Abelian Groups

Theorem 5.3.7 (The First Isomorphism Theorem). *Let $\phi : G \rightarrow H$ be a homomorphism with kernel K . Then $G/K \cong \text{im}(\phi)$. More precisely, there is an isomorphism $\bar{\phi} : G/K \rightarrow \text{im}(\phi)$ defined by $\bar{\phi}(K + g) = \phi(g)$ for all $g \in G$.*

Proof. The map $\bar{\phi}$ is a well-defined homomorphism by Lemma 5.3.6. Clearly, $\text{im}(\bar{\phi}) = \text{im}(\phi)$. Finally,

$$\bar{\phi}(K + g) = 0_H \iff \phi(g) = 0_H \iff g \in K \iff K + g = K + 0_G = 0_{G/K}.$$

By Proposition 5.3.4, $\bar{\phi}$ is injective. Thus $\bar{\phi} : G/K \rightarrow \text{im}(\phi)$ is an isomorphism. □

5.4 Free abelian groups

Definition 5.4.1. The direct sum \mathbb{Z}^n of n copies of \mathbb{Z} is known as a (finitely generated) *free abelian group* of rank n .

More generally, a finitely generated abelian group is called *free abelian* if it is isomorphic to \mathbb{Z}^n for some $n \geq 0$.

(The free abelian group \mathbb{Z}^0 of rank 0 is defined to be the trivial group $\{0\}$ containing the single element 0.)

The groups \mathbb{Z}^n have many properties in common with vector spaces such as \mathbb{R}^n , but we must expect some differences, because \mathbb{Z} is not a field.

Given we wish to exploit this connection we choose to write elements of \mathbb{Z}^n as columns, rather than rows. So $\begin{pmatrix} 1 \\ 0 \end{pmatrix} \in \mathbb{Z}^2$ etc.

We then define the standard basis of \mathbb{Z}^n exactly as for \mathbb{R}^n ; that is, $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, where \mathbf{x}_i has a 1 in its i -th component and 0 in the other components. This has the same properties as a basis of a vector space; i.e. it is linearly independent and spans \mathbb{Z}^n .

Definition 5.4.2. Elements x_1, \dots, x_n of an abelian group G are called *linearly independent* if $\alpha_1 x_1 + \dots + \alpha_n x_n = 0_G$ for $\alpha_1, \dots, \alpha_n \in \mathbb{Z}$ implies that $\alpha_1 = \alpha_2 = \dots = \alpha_n = 0_{\mathbb{Z}}$.

Definition 5.4.3. Elements x_1, \dots, x_n form a *free basis* of the abelian group G if and only if they are linearly independent and generate (span) G .

Example. It's clear that the standard basis $\mathbf{x}_1 = (1, 0, \dots, 0)^T, \mathbf{x}_2 = (0, 1, \dots, 0)^T, \dots, \mathbf{x}_n = (0, 0, \dots, 1)^T$ is indeed a free basis of \mathbb{Z}^n but there are others; for instance, $\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\}$ is a free basis of \mathbb{Z}^2 .

It's important to notice, though, that a subset of \mathbb{Z}^n which is a basis of \mathbb{Q}^n need not be a free basis of \mathbb{Z}^n . For instance, $\left\{ \begin{pmatrix} 2 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 2 \end{pmatrix} \right\}$ is not a free basis of \mathbb{Z}^2 , since we can't write all elements of \mathbb{Z}^2 as linear combinations of these vectors with *integer* coefficients – we'll need to divide by 2 at some point. This also shows that a set of n linearly independent elements of \mathbb{Z}^n needn't be a free basis.

5 Finitely Generated Abelian Groups

Now consider elements g_1, \dots, g_n of an abelian group G . It is possible to extend the assignment $\phi(\mathbf{x}_i) = g_i$ to a group homomorphism $\phi : \mathbb{Z}^n \rightarrow G$. As a function we define $\phi((a_1, a_2, \dots, a_n)^T) = \sum_{i=1}^n a_i g_i$. We leave the proof of the following result as an exercise.

Proposition 5.4.4. (i) *The function ϕ is a group homomorphism.*

(ii) *The set of elements $\{g_i\}$ are linearly independent if and only if ϕ is injective.*

(iii) *The set of elements $\{g_i\}$ span G if and only if ϕ is surjective.*

(iv) *The set of elements $\{g_i\}$ form a free basis of G if and only if ϕ is an isomorphism.*

Note that this proposition makes perfect sense for vector spaces. Also note that the last statement implies that g_1, \dots, g_n is a free basis of G if and only if every element $g \in G$ has a unique expression $g = \alpha_1 g_1 + \dots + \alpha_n g_n$ with $\alpha_i \in \mathbb{Z}$, very much like for vector spaces.

Before Proposition 5.4.4 we were trying to extend the assignment $\phi(\mathbf{x}_i) = g_i$ to a group homomorphism $\phi : \mathbb{Z}^n \rightarrow G$. Note that the extension we wrote is unique. This is the key to the next corollary. The details of the proof are left to the reader.

Corollary 5.4.5 (Universal property of the free abelian group). *Let G be a free abelian group with a free basis g_1, \dots, g_n . Let H be an abelian group and $a_1, \dots, a_n \in H$. Then there exists a unique group homomorphism $\phi : G \rightarrow H$ such that $\phi(g_i) = a_i$ for all i .*

As for finite dimensional vector spaces, it turns out that any two free bases of a free abelian group have the same size, but this has to be proved. It will follow directly from the next theorem.

Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ be the standard free basis of \mathbb{Z}^n , and let $\mathbf{y}_1, \dots, \mathbf{y}_m$ be another free basis. As in Linear Algebra, we can define the associated change of basis matrix P (with original basis $\{\mathbf{x}_i\}$ and new basis $\{\mathbf{y}_i\}$), where the columns of P are \mathbf{y}_i ; that is, they express \mathbf{y}_i in terms of \mathbf{x}_i . For example, if $n = m = 2$, $\mathbf{y}_1 = \begin{pmatrix} 2 \\ 7 \end{pmatrix}$, $\mathbf{y}_2 = \begin{pmatrix} 1 \\ 4 \end{pmatrix}$, then $P = \begin{pmatrix} 2 & 1 \\ 7 & 4 \end{pmatrix}$. In general, $P = (\rho_{ij})$ is an $n \times m$ matrix with $\mathbf{y}_j = \sum_{i=1}^n \rho_{ij} \mathbf{x}_i$ for $1 \leq j \leq m$.

Theorem 5.4.6. *Let $\mathbf{y}_1, \dots, \mathbf{y}_m \in \mathbb{Z}^n$ with $\mathbf{y}_j = \sum_{i=1}^n \rho_{ij} \mathbf{x}_i$ for $1 \leq j \leq m$. Then the following are equivalent:*

(i) $\mathbf{y}_1, \dots, \mathbf{y}_m$ is a free basis of \mathbb{Z}^n ;

(ii) $n = m$ and P is an invertible matrix such that P^{-1} has entries in \mathbb{Z} ;

(iii) $n = m$ and $\det(P) = \pm 1$.

(A matrix $P \in \mathbb{Z}^{n,n}$ with $\det(P) = \pm 1$ is called *unimodular*.)

Proof. (i) \Rightarrow (ii). If $\mathbf{y}_1, \dots, \mathbf{y}_m$ is a free basis of \mathbb{Z}^n then it spans \mathbb{Z}^n , so there is an $m \times n$ matrix $T = (\tau_{ij})$ with $\mathbf{x}_k = \sum_{j=1}^m \tau_{jk} \mathbf{y}_j$ for $1 \leq k \leq n$. Hence

$$\mathbf{x}_k = \sum_{j=1}^m \tau_{jk} \mathbf{y}_j = \sum_{j=1}^m \tau_{jk} \sum_{i=1}^n \rho_{ij} \mathbf{x}_i = \sum_{i=1}^n \left(\sum_{j=1}^m \rho_{ij} \tau_{jk} \right) \mathbf{x}_i,$$

5 Finitely Generated Abelian Groups

and, since $\mathbf{x}_1, \dots, \mathbf{x}_n$ is a free basis, this implies that $\sum_{j=1}^m \rho_{ij} \tau_{jk} = 1$ when $i = k$ and 0 when $i \neq k$. In other words $PT = I_n$, and similarly $TP = I_m$, so P and T are inverse matrices. But we can think of P and T as inverse matrices over the field \mathbb{Q} , so it follows from First Year Linear Algebra that $m = n$, and $T = P^{-1}$ has entries in \mathbb{Z} .

(ii) \Rightarrow (i). If $T = P^{-1}$ has entries in \mathbb{Z} then, again thinking of them as matrices over the field \mathbb{Q} , $\text{rank}(P) = n$, so the columns of P are linearly independent over \mathbb{Q} and hence also over \mathbb{Z} . Since the columns of P are just the column vectors representing $\mathbf{y}_1, \dots, \mathbf{y}_m$, this tells us that $\mathbf{y}_1, \dots, \mathbf{y}_m$ are linearly independent.

Using $PT = I_n$, for $1 \leq k \leq n$ we have

$$\sum_{j=1}^m \tau_{jk} \mathbf{y}_j = \sum_{j=1}^m \tau_{jk} \sum_{i=1}^n \rho_{ij} \mathbf{x}_i = \sum_{i=1}^n \left(\sum_{j=1}^m \rho_{ij} \tau_{jk} \right) \mathbf{x}_i = \mathbf{x}_k,$$

because $\sum_{j=1}^m \rho_{ij} \tau_{jk}$ is equal to 1 when $i = k$ and 0 when $i \neq k$. Since $\mathbf{x}_1, \dots, \mathbf{x}_n$ spans \mathbb{Z}^n , and we can express each \mathbf{x}_k as a linear combination of $\mathbf{y}_1, \dots, \mathbf{y}_m$, it follows that $\mathbf{y}_1, \dots, \mathbf{y}_m$ span \mathbb{Z}^n and hence form a free basis of \mathbb{Z}^n .

(ii) \Rightarrow (iii). If $T = P^{-1}$ has entries in \mathbb{Z} , then $\det(PT) = \det(P) \det(T) = \det(I_n) = 1$, and since $\det(P), \det(T) \in \mathbb{Z}$, this implies $\det(P) = \pm 1$.

(iii) \Rightarrow (ii). From First year Linear Algebra, $P^{-1} = \frac{1}{\det(P)} \text{adj}(P)$, so $\det(P) = \pm 1$ implies that P^{-1} has entries in \mathbb{Z} . □

Example. If $n = 2$ and $\mathbf{y}_1 = \begin{pmatrix} 2 \\ 7 \end{pmatrix}$, $\mathbf{y}_2 = \begin{pmatrix} 1 \\ 4 \end{pmatrix}$, then $\det(P) = 8 - 7 = 1$, so $\mathbf{y}_1, \mathbf{y}_2$ is a free basis of \mathbb{Z}^2 .

But, if $\mathbf{y}_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, $\mathbf{y}_2 = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$, then $\det(P) = 2$, so $\mathbf{y}_1, \mathbf{y}_2$ is not a free basis of \mathbb{Z}^2 .

Recall that in Linear Algebra over a field, any set of n linearly independent vectors in a vector space V of dimension n form a basis of V . This example shows that this result is not true in \mathbb{Z}^n , because \mathbf{y}_1 and \mathbf{y}_2 are linearly independent but do not span \mathbb{Z}^2 .

But as in Linear Algebra, for $\mathbf{v} \in \mathbb{Z}^n$, if \mathbf{x} and \mathbf{y} are the column vectors representing \mathbf{v} using free bases $\mathbf{x}_1, \dots, \mathbf{x}_n$ and $\mathbf{y}_1, \dots, \mathbf{y}_n$, respectively, then we have $\mathbf{x} = P\mathbf{y}$, so $\mathbf{y} = P^{-1}\mathbf{x}$.

5.5 Unimodular elementary row and column operations and the unimodular Smith normal form for integer matrices

We interrupt our discussion of finitely generated abelian groups at this stage to investigate how the row and column reduction process of Linear Algebra can be adapted to matrices over \mathbb{Z} . Recall from MA106 that we can use elementary row and column operations to reduce an $m \times n$ matrix of rank r over a field K to a matrix $B = (\beta_{ij})$ with $\beta_{ii} = 1$ for $1 \leq i \leq r$ and $\beta_{ij} = 0$ otherwise. We called this the *Smith normal form* of the matrix. We can do something similar over \mathbb{Z} , but the non-zero elements β_{ii} will not necessarily all be equal to 1.

5 Finitely Generated Abelian Groups

The reason that we disallowed $\lambda = 0$ for the row and column operations (R3) and (C3) (multiply a row or column by a scalar λ) was that we wanted all of our elementary operations to be reversible. When performed over \mathbb{Z} , (R1), (C1), (R2) and (C2) are reversible, but (R3) and (C3) are reversible only when $\lambda = \pm 1$. So, if A is an $m \times n$ matrix over \mathbb{Z} , then we define the three types of *unimodular* elementary row and column operations as follows:

- (UR1): Replace some row \mathbf{r}_i of A by $\mathbf{r}_i + t\mathbf{r}_j$, where $j \neq i$ and $t \in \mathbb{Z}$;
- (UR2): Interchange two rows \mathbf{r}_i and \mathbf{r}_j of A ;
- (UR3): Replace some row \mathbf{r}_i of A by $-\mathbf{r}_i$.
- (UC1): Replace some column \mathbf{c}_i of A by $\mathbf{c}_i + t\mathbf{c}_j$, where $j \neq i$ and $t \in \mathbb{Z}$;
- (UC2): Interchange two columns \mathbf{c}_i and \mathbf{c}_j of A ;
- (UC3): Replace some column \mathbf{c}_i of A by $-\mathbf{c}_i$.

Recall from MA106 that performing elementary row or column operations on a matrix A corresponds to multiplying A on the left or right, respectively, by an elementary matrix. These elementary matrices all have determinant ± 1 (1 for (UR1) and -1 for (UR2) and (UR3)), so are unimodular matrices over \mathbb{Z} .

By checking what left and right multiplying by these elementary matrices do we can check that the unimodular row and column operations correspond to the following change of bases, where $\mathbf{e}_1, \dots, \mathbf{e}_n$ is a basis for \mathbb{Z}^n (the domain of the linear map which A represents) and $\mathbf{f}_1, \dots, \mathbf{f}_m$ is a basis for \mathbb{Z}^m (the target of the linear map). Notice that column operations correspond to changing the basis of the domain, whereas row operations correspond to changing the basis of the target.

- | | | |
|--|--|---|
| (UC1): $\mathbf{e}_i \rightarrow \mathbf{e}_i + t\mathbf{e}_j$; | (UC2): $\mathbf{e}_i \leftrightarrow \mathbf{e}_j$; | (UC3): $\mathbf{e}_i \rightarrow -\mathbf{e}_i$. |
| (UR1): $\mathbf{f}_j \rightarrow \mathbf{f}_j - t\mathbf{f}_i$; | (UR2): $\mathbf{f}_i \leftrightarrow \mathbf{f}_j$; | (UR3): $\mathbf{f}_i \rightarrow -\mathbf{f}_i$. |

Theorem 5.5.1. *Let A be an $m \times n$ matrix over \mathbb{Z} with rank r . Then, by using a sequence of unimodular elementary row and column operations, we can reduce A to a matrix $B = (\beta_{ij})$ with $\beta_{ii} = d_i$ for $1 \leq i \leq r$ and $\beta_{ij} = 0$ otherwise, and where the integers d_i satisfy $d_i > 0$ for $1 \leq i \leq r$, and $d_i | d_{i+1}$ for $1 \leq i < r$. Subject to these conditions, the d_i are uniquely determined by the matrix A .*

Proof. We shall not prove the uniqueness part here. The fact that the number of non-zero β_{ii} is the rank of A follows from the fact that unimodular row and column operations do not change the rank. We use induction on $m + n$. The base case is $m = n = 1$, where there is nothing to prove. Also if A is the zero matrix then there is nothing to prove, so assume not.

Let d be the smallest entry with $d > 0$ in any matrix $C = (\gamma_{ij})$ that we can obtain from A by using unimodular elementary row and column operations. By using (UR2) and (UC2), we can move d to position $(1, 1)$ and hence assume that $\gamma_{11} = d$. If d does not divide γ_{1j} for some $j > 0$, then we can write $\gamma_{1j} = qd + r$ with $q, r \in \mathbb{Z}$ and $0 < r < d$, and then replacing the j -th column \mathbf{c}_j of C by $\mathbf{c}_j - q\mathbf{c}_1$ results in the entry r in position $(1, j)$, contrary to the choice of d . Hence $d | \gamma_{1j}$ for $2 \leq j \leq n$ and similarly $d | \gamma_{i1}$ for $2 \leq i \leq m$.

5 Finitely Generated Abelian Groups

Now, if $\gamma_{1j} = qd$, then replacing \mathbf{c}_j of C by $\mathbf{c}_j - q\mathbf{c}_1$ results in entry 0 position $(1, j)$. So we can assume that $\gamma_{1j} = 0$ for $2 \leq j \leq n$ and $\gamma_{i1} = 0$ for $2 \leq i \leq m$. If $m = 1$ or $n = 1$, then we are done. Otherwise, we have $C = (d) \oplus C'$ for some $(m-1) \times (n-1)$ matrix C' . By inductive hypothesis, the result of the theorem applies to C' , so by applying unimodular row and column operations to C which do not involve the first row or column, we can reduce C to $D = (\delta_{ij})$, which satisfies $\delta_{11} = d$, $\delta_{ii} = d_i > 0$ for $2 \leq i \leq r$, and $\delta_{ij} = 0$ otherwise, where $d_i | d_{i+1}$ for $2 \leq i < r$. To complete the proof, we still have to show that $d | d_2$. If not, then adding row 2 to row 1 results in d_2 in position $(1,2)$ not divisible by d , and we obtain a contradiction as before. \square

Definition 5.5.2. Let A be an $m \times n$ matrix over \mathbb{Z} with rank r . The uniquely determined diagonal matrix in Theorem 5.5.1, where the d_i satisfy $d_i > 0$ for $1 \leq i \leq r$, and $d_i | d_{i+1}$ for $1 \leq i < r$ is called the *unimodular Smith normal form* of A or just SNF for short.

So how do we find the unimodular Smith normal form of a matrix then? The general strategy is to reduce the size of entries in the first row and column, until the $(1,1)$ -entry divides all other entries in the first row and column. Then we can clear all of these other entries with repeated use of (UR1) and (UC1). Let us elaborate on this strategy. First we state a useful result in finding the SNF.

Lemma 5.5.3. Let $A \in \mathbb{Z}^{m,n}$ with unimodular Smith normal form S having non-zero diagonal entries d_1, \dots, d_r . Then the greatest common divisor of all of the entries of A is equal to d_1 . (Our convention is that $\gcd(r, 0) = r$ for all integers $r \geq 1$ and recall that the greatest common divisor is always non-negative).

Proof. By our convention, the greatest common divisor of all of the entries of S is equal to d_1 (since d_1 divides d_2, \dots, d_r by definition of the SNF). So it suffices to prove that applying any of the six unimodular elementary row and column operations to a matrix $B \in \mathbb{Z}^{m,n}$ preserves the greatest common divisor of all the entries. We'll show this for the row operations, the column operations are very similar. Firstly, applying (UR2) to B does not change the set of entries in B , it only permutes them, so clearly the greatest common divisor of the entries is preserved. Similarly for (UR3), multiplying a row by ± 1 only changes the signs of some of the entries which also does not affect the GCD of all the entries of B . Finally, we consider (UR1). This takes a row r_i of B and replaces it with $r_i + tr_j$ for some $i \neq j$ and $t \in \mathbb{Z}$. It suffices for us to check that the GCD of the entries in rows r_i and r_j is the same as the GCD of the entries in $r_i + tr_j$ and r_j . And to show this is true it suffices to realise that $\gcd(b_{ik}, b_{jk}) = \gcd(b_{ik} + tb_{jk}, b_{jk})$ for all $1 \leq k \leq n$. Indeed, it is clear that $\gcd(b_{ik}, b_{jk}) | \gcd(b_{ik} + tb_{jk}, b_{jk})$ and if $g | b_{jk}$ and $g | b_{ik} + tb_{jk}$ then $g | b_{ik} + tb_{jk} - tb_{jk} = b_{ik}$ and so $\gcd(b_{ik} + tb_{jk}, b_{jk}) | \gcd(b_{ik}, b_{jk})$, as required. \square

One can generalise this to the greatest common divisor of $k \times k$ minors. You may want to think about what they tell you about the SNF. This is one way to prove uniqueness of the SNF (see Example Sheet 9). Now we present a strategy for finding the SNF of an integer matrix.

Strategy for finding the SNF of $A \in \mathbb{Z}^{m,n}$

Step 1: Find d_1 , using Lemma 5.5.3.

5 Finitely Generated Abelian Groups

Step 2: If d_1 occurs as an entry in A move it to the $(1, 1)$ entry using (UC2) and (UR2). If $-d_1$ occurs then use (UR3) or (UC3) and then move it into the $(1, 1)$ entry. This is the easier scenario.

If d_1 does not occur then we need to do some (or lots) of division with remainder. Let x be the smallest entry (with respect to absolute value) occurring in A , say in position (i, j) .

Now, we again have a dichotomy into an easier case and a harder case. Does x divide everything else in the i th row and j th column? If not, let y be an entry that x does not divide, in position (k, l) with $k = i$ or $l = j$. Then $y = sx + r$ with $r < x$ and using (UC1) or (UR1) we can obtain $r = y - sx$ as an entry of A (add $-sx$ of row i to row k or $-sx$ of column j to column l). At this point we have reduced the smallest entry in \tilde{A} (the matrix obtained from A by doing the unimodular row or column operation required) and so we can start Step 2 again.

Now it remains to deal with the case where x divides everything else in row i and column j . We start by clearing all of these entries to 0. This is straightforward using (UC1) and (UR1), and is possible since x divides each entry:

$$r_d \rightarrow r_d - \frac{a_{d,j}}{x} r_i, \quad d \neq i$$

and

$$c_d \rightarrow c_d - \frac{a_{i,d}}{x} c_j, \quad d \neq j$$

do the job. We know there still exists an element in \tilde{A} which is not divisible by x , again let y be such an entry, in position (k, l) . Then looking at the intersection of row i and k with column j and l we find a 2×2 -matrix which looks as follows (assuming $i < k$ and $j < l$, otherwise the picture is similar but x and y are in different positions, still opposite each other diagonally):

$$\begin{pmatrix} x & 0 \\ 0 & y \end{pmatrix}.$$

We want $r = y - sx$ with $r < x$ to appear in \tilde{A} . We can do this in two moves. First use (UR1) to add $-s$ times row i to row k . We then use (UC1) to add column j to column l . This will result in r appearing where y was before:

$$\begin{pmatrix} x & -sx \\ -sx & r = y - sx \end{pmatrix}.$$

(There are other ways to do this, feel free to experiment!) Again, we have reduced the smallest entry in \tilde{A} and so we can start Step 2 again.

Step 3: With d_1 in the $(1, 1)$ entry, we can use it to clear everything else in the first row and column since d_1 divides all entries in the matrix (just as we did in Step 2).

Step 4: The matrix \tilde{A} now has entry d_1 in the $(1, 1)$ entry and 0s elsewhere in row and column 1. We now go back to Step 1, working on the $m - 1 \times n - 1$ matrix in the bottom right hand corner of \tilde{A} . We can repeat Steps 1 to 3 without changing the entries in row 1 and column 1. Therefore, repeating this process will terminate and will yield the SNF of A .

Example 18. $A = \begin{pmatrix} 42 & 21 \\ -35 & -14 \end{pmatrix}$.

5 Finitely Generated Abelian Groups

We calculate that $d_1 = \gcd(42, 21, -35, -14) = 7$. It does not appear in A so we need to use unimodular row and column operations to make that happen. It is often easy in practice to see how to do this and the general strategy should be seen as a guide rather than an algorithm you must follow. On this occasion we notice that dividing -35 by -14 we get remainder $-7 = -d_1$. We can then negate row 2 and we see 7 occurring in the matrix. Everything is straightforward from that point onwards.

Matrix	Operation	Matrix	Operation
$\begin{pmatrix} 42 & 21 \\ -35 & -14 \end{pmatrix}$	$\mathbf{c}_1 \rightarrow \mathbf{c}_1 - 2\mathbf{c}_2$	$\begin{pmatrix} 0 & 21 \\ -7 & -14 \end{pmatrix}$	$\mathbf{r}_2 \rightarrow -\mathbf{r}_2$ $\mathbf{r}_1 \leftrightarrow \mathbf{r}_2$
$\begin{pmatrix} 7 & 14 \\ 0 & 21 \end{pmatrix}$	$\mathbf{c}_2 \rightarrow \mathbf{c}_2 - 2\mathbf{c}_1$	$\begin{pmatrix} 7 & 0 \\ 0 & 21 \end{pmatrix}$	

Example 19. $A = \begin{pmatrix} -18 & -18 & -18 & 90 \\ 54 & 12 & 45 & 48 \\ 9 & -6 & 6 & 63 \\ 18 & 6 & 15 & 12 \end{pmatrix}$.

This time we want to notice what the greatest common divisor of the entries is without doing too many calculations. Firstly, we spot 9 and 6, so d_1 must be 1 or 3. Looking at all the other entries we see they are all divisible by 3. So $d_1 = 3$. Again, this does not appear in A so we need to use unimodular row and column operations to make that happen. We spot 9 and 6 in row 3, so we can get 3 to appear just by adding the negative of column 3 to column 1 (plenty of other choices here!). We carry on in this way to obtain the SNF.

Matrix	Operation	Matrix	Operation
$\begin{pmatrix} -18 & -18 & -18 & 90 \\ 54 & 12 & 45 & 48 \\ 9 & -6 & 6 & 63 \\ 18 & 6 & 15 & 12 \end{pmatrix}$	$\mathbf{c}_1 \rightarrow \mathbf{c}_1 - \mathbf{c}_3$	$\begin{pmatrix} 0 & -18 & -18 & 90 \\ 9 & 12 & 45 & 48 \\ 3 & -6 & 6 & 63 \\ 3 & 6 & 15 & 12 \end{pmatrix}$	$\mathbf{r}_1 \leftrightarrow \mathbf{r}_4$
$\begin{pmatrix} 3 & 6 & 15 & 12 \\ 9 & 12 & 45 & 48 \\ 3 & -6 & 6 & 63 \\ 0 & -18 & -18 & 90 \end{pmatrix}$	$\mathbf{r}_2 \rightarrow \mathbf{r}_2 - 3\mathbf{r}_1$ $\mathbf{r}_3 \rightarrow \mathbf{r}_3 - \mathbf{r}_1$	$\begin{pmatrix} 3 & 6 & 15 & 12 \\ 0 & -6 & 0 & 12 \\ 0 & -12 & -9 & 51 \\ 0 & -18 & -18 & 90 \end{pmatrix}$	$\mathbf{c}_2 \rightarrow \mathbf{c}_2 - 2\mathbf{c}_1$ $\mathbf{c}_3 \rightarrow \mathbf{c}_3 - 5\mathbf{c}_1$ $\mathbf{c}_4 \rightarrow \mathbf{c}_4 - 4\mathbf{c}_1$
$\begin{pmatrix} 3 & 0 & 0 & 0 \\ 0 & -6 & 0 & 12 \\ 0 & -12 & -9 & 51 \\ 0 & -18 & -18 & 90 \end{pmatrix}$	$\mathbf{c}_2 \rightarrow -\mathbf{c}_2$ $\mathbf{c}_2 \rightarrow \mathbf{c}_2 + \mathbf{c}_3$	$\begin{pmatrix} 3 & 0 & 0 & 0 \\ 0 & 6 & 0 & 12 \\ 0 & 3 & -9 & 51 \\ 0 & 0 & -18 & 90 \end{pmatrix}$	$\mathbf{r}_2 \leftrightarrow \mathbf{r}_3$
$\begin{pmatrix} 3 & 0 & 0 & 0 \\ 0 & 3 & -9 & 51 \\ 0 & 6 & 0 & 12 \\ 0 & 0 & -18 & 90 \end{pmatrix}$	$\mathbf{r}_3 \rightarrow \mathbf{r}_3 - 2\mathbf{r}_2$	$\begin{pmatrix} 3 & 0 & 0 & 0 \\ 0 & 3 & -9 & 51 \\ 0 & 0 & 18 & -90 \\ 0 & 0 & -18 & 90 \end{pmatrix}$	$\mathbf{c}_3 \rightarrow \mathbf{c}_3 + 3\mathbf{c}_2$ $\mathbf{c}_4 \rightarrow \mathbf{c}_4 - 17\mathbf{c}_2$

5 Finitely Generated Abelian Groups

$$\begin{pmatrix} 3 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 18 & -90 \\ 0 & 0 & -18 & 90 \end{pmatrix} \quad \begin{array}{l} \mathbf{c}_4 \rightarrow \mathbf{c}_4 + 5\mathbf{c}_3 \\ \mathbf{r}_4 \rightarrow \mathbf{r}_4 + \mathbf{r}_3 \end{array} \quad \begin{pmatrix} 3 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 18 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

Note: There is also a generalisation to integer matrices of the the row reduced normal form from Linear Algebra, where only row operations are allowed. This is known as the *Hermite Normal Form* and is more complicated.

5.6 Subgroups of free abelian groups

Proposition 5.6.1. *Any subgroup of a finitely generated abelian group is finitely generated.*

Proof. Let $K < G$ with G an abelian group generated by x_1, \dots, x_n . We shall prove by induction on n that K can be generated by at most n elements. If $n = 1$ then G is cyclic. Write $G = \{sx \mid s \in \mathbb{Z}\}$. Let m be the smallest positive number such that $mx \in K$. If such a number does not exist then $K = \{0\}$. Otherwise, $K \supseteq \{smx \mid s \in \mathbb{Z}\}$. The opposite inclusion follows using division with a remainder: write $t = qm + r$ with $0 \leq r < m$. Then $tx \in K$ if and only if $rx = (t - mq)x \in K$ if and only if $r = 0$ due to minimality of m . In both cases K is cyclic.

Suppose $n > 1$, and let H be the subgroup of G generated by x_1, \dots, x_{n-1} . By induction, $K \cap H$ is generated by y_1, \dots, y_{m-1} , say, with $m \leq n$. If $K \leq H$, then $K = K \cap H$ and we are done, so suppose not.

Then there exist elements of the form $h + tx_n \in K$ with $h \in H$ and $t \neq 0$. Since $-(h + tx_n) \in K$, we can assume that $t > 0$. Choose such an element $y_m = h + tx_n \in K$ with t minimal subject to $t > 0$. We claim that K is generated by y_1, \dots, y_m , which will complete the proof. Let $k \in K$. Then $k = h' + ux_n$ with $h' \in H$ and $u \in \mathbb{Z}$. If t does not divide u then we can write $u = tq + r$ with $q, r \in \mathbb{Z}$ and $0 < r < t$, and then $k - qy_m = (h' - qh) + rx_n \in K$, contrary to the choice of t . So $t \mid u$ and hence $u = tq$ and $k - qy_m \in K \cap H$. But $K \cap H$ is generated by y_1, \dots, y_{m-1} , so we are done. \square

Now let H be a subgroup of the free abelian group \mathbb{Z}^n , and suppose that H is generated by $\mathbf{v}_1, \dots, \mathbf{v}_m$. Then H can be represented by an $n \times m$ matrix A in which the columns are $\mathbf{v}_1, \dots, \mathbf{v}_m$.

Example 20. If $n = 3$ and H is generated by $\mathbf{v}_1 = \begin{pmatrix} 1 \\ 3 \\ -1 \end{pmatrix}$ and $\mathbf{v}_2 = \begin{pmatrix} 2 \\ 0 \\ 1 \end{pmatrix}$, then

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 0 \\ -1 & 1 \end{pmatrix}.$$

As we saw above, if we use a different free basis $\mathbf{y}_1, \dots, \mathbf{y}_n$ of \mathbb{Z}^n with basis change matrix P , then each column \mathbf{v}_j of A is replaced by $P^{-1}\mathbf{v}_j$, and hence A itself is replaced by $P^{-1}A$.

5 Finitely Generated Abelian Groups

So in Example 20, if we use the basis $\mathbf{y}_1 = \begin{pmatrix} 0 \\ -1 \\ 0 \end{pmatrix}$, $\mathbf{y}_2 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$, $\mathbf{y}_3 = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}$ of \mathbb{Z}^3 , then

$$P = \begin{pmatrix} 0 & 1 & 1 \\ -1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}, \quad P^{-1} = \begin{pmatrix} 1 & -1 & -1 \\ 0 & 0 & 1 \\ 1 & 0 & -1 \end{pmatrix}, \quad P^{-1}A = \begin{pmatrix} -1 & 1 \\ -1 & 1 \\ 2 & 1 \end{pmatrix}.$$

For example, the first column $\begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix}$ of $P^{-1}A$ represents $-\mathbf{y}_1 - \mathbf{y}_2 + 2\mathbf{y}_3 = \begin{pmatrix} 1 \\ 3 \\ -1 \end{pmatrix} = \mathbf{v}_1$.

In particular, if we perform a unimodular elementary row operation on A , then the resulting matrix represents the same subgroup H of \mathbb{Z}^n but using a different free basis of \mathbb{Z}^n .

We can clearly replace a generator \mathbf{v}_i of H by $\mathbf{v}_i + r\mathbf{v}_j$ for $r \in \mathbb{Z}$ without changing the subgroup H that is generated. We can also interchange two of the generators or replace one of the generators \mathbf{v}_i by $-\mathbf{v}_i$ without changing H . In other words, performing a unimodular elementary column operation on A amounts to changing the generating set for H , so again the resulting matrix still represents the same subgroup H of \mathbb{Z}^n .

Summing up, we have:

Proposition 5.6.2. *Suppose that the subgroup H of \mathbb{Z}^n is represented by the matrix $A \in \mathbb{Z}^{n,m}$. Then if the matrix $B \in \mathbb{Z}^{n,m}$ is obtained by performing a sequence of unimodular row and column operations on A , then B represents the same subgroup H of \mathbb{Z}^n using a (possibly) different free basis of \mathbb{Z}^n .*

In particular, by Theorem 5.5.1, we can transform A to its unimodular Smith normal form B . So, then if B represents H with the free basis $\mathbf{y}_1, \dots, \mathbf{y}_n$ of \mathbb{Z}^n , then the r non-zero columns of B correspond to the elements $d_1\mathbf{y}_1, d_2\mathbf{y}_2, \dots, d_r\mathbf{y}_r$ of \mathbb{Z}^n . So we have:

Theorem 5.6.3. *Let H be a subgroup of \mathbb{Z}^n . Then there exists a free basis $\mathbf{y}_1, \dots, \mathbf{y}_n$ of \mathbb{Z}^n such that $H = \langle d_1\mathbf{y}_1, d_2\mathbf{y}_2, \dots, d_r\mathbf{y}_r \rangle$, where each $d_i > 0$ and $d_i | d_{i+1}$ for $1 \leq i < r$.*

In Example 20, it is straightforward to calculate the unimodular Smith normal form of A , which is $\begin{pmatrix} 1 & 0 \\ 0 & 3 \\ 0 & 0 \end{pmatrix}$, so $H = \langle \mathbf{y}_1, 3\mathbf{y}_2 \rangle$.

By keeping track of the unimodular row operations carried out, we can, if we need to, find the free basis $\mathbf{y}_1, \dots, \mathbf{y}_n$ of \mathbb{Z}^n such that H has this nice form. Using the formulae in Section 5.5, noting that we start from the standard free basis, we can do this in Example 20.

Matrix	Operation	New free basis
$\begin{pmatrix} 1 & 2 \\ 3 & 0 \\ -1 & 1 \end{pmatrix}$	$\mathbf{r}_2 \rightarrow \mathbf{r}_2 - 3\mathbf{r}_1$	$\mathbf{y}_1 = \begin{pmatrix} 1 \\ 3 \\ 0 \end{pmatrix}, \mathbf{y}_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \mathbf{y}_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$

5 Finitely Generated Abelian Groups

$$\begin{array}{l}
 \begin{pmatrix} 1 & 2 \\ 0 & -6 \\ -1 & 1 \end{pmatrix} \quad \mathbf{r}_3 \rightarrow \mathbf{r}_3 + \mathbf{r}_1 \quad \mathbf{y}_1 = \begin{pmatrix} 1 \\ 3 \\ -1 \end{pmatrix}, \mathbf{y}_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \mathbf{y}_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \\
 \\
 \begin{pmatrix} 1 & 2 \\ 0 & -6 \\ 0 & 3 \end{pmatrix} \quad \mathbf{c}_2 \rightarrow \mathbf{c}_2 - 2\mathbf{c}_1 \quad \mathbf{y}_1 = \begin{pmatrix} 1 \\ 3 \\ -1 \end{pmatrix}, \mathbf{y}_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \mathbf{y}_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \\
 \\
 \begin{pmatrix} 1 & 0 \\ 0 & -6 \\ 0 & 3 \end{pmatrix} \quad \mathbf{r}_2 \leftrightarrow \mathbf{r}_3 \quad \mathbf{y}_1 = \begin{pmatrix} 1 \\ 3 \\ -1 \end{pmatrix}, \mathbf{y}_2 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}, \mathbf{y}_3 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \\
 \\
 \begin{pmatrix} 1 & 0 \\ 0 & 3 \\ 0 & -6 \end{pmatrix} \quad \mathbf{r}_3 \rightarrow \mathbf{r}_3 + 2\mathbf{r}_2 \quad \mathbf{y}_1 = \begin{pmatrix} 1 \\ 3 \\ -1 \end{pmatrix}, \mathbf{y}_2 = \begin{pmatrix} 0 \\ -2 \\ 1 \end{pmatrix}, \mathbf{y}_3 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \\
 \\
 \begin{pmatrix} 1 & 0 \\ 0 & 3 \\ 0 & 0 \end{pmatrix}
 \end{array}$$

5.7 General finitely generated abelian groups

Let G be a finitely generated abelian group. If G has n generators, Proposition 5.4.4 gives a surjective homomorphism $\phi : \mathbb{Z}^n \rightarrow G$. From the First isomorphism Theorem (Theorem 5.3.7) we deduce that $G \cong \mathbb{Z}^n / K$, where $K = \ker(\phi)$. So we have proved that every finitely generated abelian group is isomorphic to a quotient group of a free abelian group.

From the definition of ϕ , we see that

$$K = \{ (\alpha_1, \alpha_2, \dots, \alpha_n)^T \in \mathbb{Z}^n \mid \alpha_1 x_1 + \dots + \alpha_n x_n = 0_G \}.$$

By Theorem 5.6.1, this subgroup K is generated by finitely many elements $\mathbf{v}_1, \dots, \mathbf{v}_m$ of \mathbb{Z}^n . The notation

$$\langle \mathbf{x}_1, \dots, \mathbf{x}_n \mid \mathbf{v}_1, \dots, \mathbf{v}_m \rangle$$

is often used to denote the quotient group \mathbb{Z}^n / K , so we have

$$G \cong \langle \mathbf{x}_1, \dots, \mathbf{x}_n \mid \mathbf{v}_1, \dots, \mathbf{v}_m \rangle.$$

Now we can apply Theorem 5.6.3 to this subgroup K , and deduce that there is a free basis $\mathbf{y}_1, \dots, \mathbf{y}_n$ of \mathbb{Z}^n such that $K = \langle d_1 \mathbf{y}_1, \dots, d_r \mathbf{y}_r \rangle$ for some $r \leq n$, where each $d_i > 0$ and $d_i \mid d_{i+1}$ for $1 \leq i < r$.

So we also have

$$G \cong \langle \mathbf{y}_1, \dots, \mathbf{y}_n \mid d_1 \mathbf{y}_1, \dots, d_r \mathbf{y}_r \rangle,$$

and G has generators y_1, \dots, y_n with $d_i y_i = 0$ for $1 \leq i \leq r$.

5 Finitely Generated Abelian Groups

Proposition 5.7.1. *The group*

$$\langle \mathbf{y}_1, \dots, \mathbf{y}_n \mid d_1 \mathbf{y}_1, \dots, d_r \mathbf{y}_r \rangle$$

is isomorphic to the direct sum of cyclic groups

$$\mathbb{Z}_{d_1} \oplus \mathbb{Z}_{d_2} \oplus \dots \oplus \mathbb{Z}_{d_r} \oplus \mathbb{Z}^{n-r}.$$

Proof. This is another application of the First Isomorphism Theorem. Let $H = \mathbb{Z}_{d_1} \oplus \mathbb{Z}_{d_2} \oplus \dots \oplus \mathbb{Z}_{d_r} \oplus \mathbb{Z}^{n-r}$, so H is generated by y_1, \dots, y_n , with $y_1 = (1, 0, \dots, 0), \dots, y_n = (0, 0, \dots, 1)$. Let us consider \mathbb{Z}^n such that $\mathbf{y}_1, \dots, \mathbf{y}_n$ is its standard free basis. Then, by Proposition 5.4.4, there is a surjective homomorphism ϕ from \mathbb{Z}^n to H for which

$$\phi(\alpha_1 \mathbf{y}_1 + \dots + \alpha_n \mathbf{y}_n) = \alpha_1 y_1 + \dots + \alpha_n y_n$$

for all $\alpha_1, \dots, \alpha_n \in \mathbb{Z}$. Then, by Theorem 5.3.7, we have $H \cong \mathbb{Z}^n / K$, with

$$K = \{ (\alpha_1, \alpha_2, \dots, \alpha_n)^T \in \mathbb{Z}^n \mid \alpha_1 y_1 + \dots + \alpha_n y_n = 0_H \}.$$

Now $\alpha_1 y_1 + \dots + \alpha_n y_n$ is the element $(\alpha_1, \alpha_2, \dots, \alpha_n)$ of H , which is the zero element if and only if α_i is the zero element of \mathbb{Z}_{d_i} for $1 \leq i \leq r$ and $\alpha_i = 0$ for $r + 1 \leq i \leq n$.

But α_i is the zero element of \mathbb{Z}_{d_i} if and only if $d_i \mid \alpha_i$, so we have

$$K = \{ (\alpha_1, \alpha_2, \dots, \alpha_r, 0, \dots, 0)^T \in \mathbb{Z}^n \mid d_i \mid \alpha_i \text{ for } 1 \leq i \leq r \}$$

which is generated by the elements $d_1 \mathbf{y}_1, \dots, d_r \mathbf{y}_r$. So

$$H \cong \mathbb{Z}^n / K = \langle \mathbf{y}_1, \dots, \mathbf{y}_n \mid d_1 \mathbf{y}_1, \dots, d_r \mathbf{y}_r \rangle.$$

□

Putting all of these results together, we get the main theorem:

Theorem 5.7.2 (The fundamental theorem of finitely generated abelian groups). *If G is a finitely generated abelian group, then G is isomorphic to a direct sum of cyclic groups. More precisely, if G is generated by n elements then, for some r with $0 \leq r \leq n$, there are integers d_1, \dots, d_r with $d_i > 0$ and $d_i \mid d_{i+1}$ such that*

$$G \cong \mathbb{Z}_{d_1} \oplus \mathbb{Z}_{d_2} \oplus \dots \oplus \mathbb{Z}_{d_r} \oplus \mathbb{Z}^{n-r}.$$

So G is isomorphic to a direct sum of r finite cyclic groups of orders d_1, \dots, d_r , and $n - r$ infinite cyclic groups.

There may be some factors \mathbb{Z}_1 , the trivial group of order 1. These can be omitted from the direct sum (except in the case when $G \cong \mathbb{Z}_1$ is trivial). It can be deduced from the uniqueness part of Theorem 5.5.1, which we did not prove, that the numbers in the sequence d_1, d_2, \dots, d_r that are greater than 1 are uniquely determined by G .

Note that, $n - r$ may be 0, which is the case if and only if G is finite. At the other extreme, if all $d_i = 1$, then G is free abelian.

5 Finitely Generated Abelian Groups

The group G corresponding to Example 18 in Section 5.5 is

$$\langle \mathbf{x}_1, \mathbf{x}_2 \mid 42\mathbf{x}_1 - 35\mathbf{x}_2, 21\mathbf{x}_1 - 14\mathbf{x}_2 \rangle$$

and we have $G \cong \mathbb{Z}_7 \oplus \mathbb{Z}_{21}$, a group of order $7 \times 21 = 147$.

The group defined by Example 19 in Section 5.5 is

$$\langle \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4 \mid \begin{array}{ll} -18\mathbf{x}_1 + 54\mathbf{x}_2 + 9\mathbf{x}_3 + 18\mathbf{x}_4, & -18\mathbf{x}_1 + 12\mathbf{x}_2 - 6\mathbf{x}_3 + 6\mathbf{x}_4, \\ -18\mathbf{x}_1 + 45\mathbf{x}_2 + 6\mathbf{x}_3 + 15\mathbf{x}_4, & 90\mathbf{x}_1 + 48\mathbf{x}_2 + 63\mathbf{x}_3 + 12\mathbf{x}_4 \end{array} \rangle,$$

which is isomorphic to $\mathbb{Z}_3 \oplus \mathbb{Z}_3 \oplus \mathbb{Z}_{18} \oplus \mathbb{Z}$, and is an infinite group with a (maximal) finite subgroup of order $3 \times 3 \times 18 = 162$,

The group defined by Example 20 in Section 5.6 is

$$\langle \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3 \mid \mathbf{x}_1 + 3\mathbf{x}_2 - \mathbf{x}_3, 2\mathbf{x}_1 + \mathbf{x}_3 \rangle,$$

and is isomorphic to $\mathbb{Z}_1 \oplus \mathbb{Z}_3 \oplus \mathbb{Z} \cong \mathbb{Z}_3 \oplus \mathbb{Z}$, so it is infinite, with a finite subgroup of order 3.

5.8 Finite abelian groups

In particular, for any finite abelian group G , we have $G \cong \mathbb{Z}_{d_1} \oplus \mathbb{Z}_{d_2} \oplus \cdots \oplus \mathbb{Z}_{d_r}$, where $d_i \mid d_{i+1}$ for $1 \leq i < r$, and $|G| = d_1 d_2 \cdots d_r$.

From the uniqueness part of Theorem 5.5.1 (which we did not prove), it follows that, if $d_i \mid d_{i+1}$ for $1 \leq i < r$ and $e_i \mid e_{i+1}$ for $1 \leq i < s$. then $\mathbb{Z}_{d_1} \oplus \mathbb{Z}_{d_2} \oplus \mathbb{Z}_{d_r} \cong \mathbb{Z}_{e_1} \oplus \mathbb{Z}_{e_2} \oplus \mathbb{Z}_{e_s}$ if and only if $r = s$ and $d_i = e_i$ for $1 \leq i \leq r$.

So the isomorphism classes of finite abelian groups of order $n > 0$ are in one-one correspondence with expressions $n = d_1 d_2 \cdots d_r$ for which $d_i \mid d_{i+1}$ for $1 \leq i < r$. This enables us to classify isomorphism classes of finite abelian groups.

Examples. 1. $n = 4$. The decompositions are 4 and 2×2 , so $G \cong \mathbb{Z}_4$ or $\mathbb{Z}_2 \oplus \mathbb{Z}_2$.

2. $n = 15$. The only decomposition is 15, so $G \cong \mathbb{Z}_{15}$ is necessarily cyclic.

3. $n = 36$. Decompositions are 36, 2×18 , 3×12 and 6×6 , so $G \cong \mathbb{Z}_{36}, \mathbb{Z}_2 \oplus \mathbb{Z}_{18}, \mathbb{Z}_3 \oplus \mathbb{Z}_{12}$ and $\mathbb{Z}_6 \oplus \mathbb{Z}_6$.

Although we have not proved in general that groups of the same order but with different decompositions of the type above are not isomorphic, this can always be done in specific examples by looking at the orders of elements.

You were asked to prove in an exercise that if $\phi : G \rightarrow H$ is an isomorphism then $|g| = |\phi(g)|$ for all $g \in G$. So isomorphic groups have the same number of elements of each order. And the following lemma gives us the tool we need to count the elements of a given order in direct sums of groups.

5 Finitely Generated Abelian Groups

Lemma 5.8.1. *Let $G = G_1 \oplus \cdots \oplus G_n$ be a finite abelian group. The order of $g = (g_1, g_2, \dots, g_n)$ is the least common multiple of the orders $|g_i|$ of the components of g .*

Proof. To start with we let $l = \text{lcm}(|g_1|, \dots, |g_n|)$ and note that $lg = (lg_1, \dots, lg_n) = (0, \dots, 0)$, since $lg_i = 0$ for all i by definition of l . This shows that $|g|$ divides l . Now suppose that $rg = 0$. Then $(rg_1, \dots, rg_n) = (0, \dots, 0)$ and so considering each component we see that r is divisible by $|g_1|$, and by $|g_2|$, etc. But that means that r is divisible by $\text{lcm}(|g_1|, \dots, |g_n|)$ and so $\text{lcm}(|g_1|, \dots, |g_n|)$ divides $|g|$, and we are done. \square

So, let's go back to the four groups of order 36 coming from the decompositions of 36 above, $G_1 = \mathbb{Z}_{36}$, $G_2 = \mathbb{Z}_2 \oplus \mathbb{Z}_{18}$, $G_3 = \mathbb{Z}_3 \oplus \mathbb{Z}_{12}$ and $G_4 = \mathbb{Z}_6 \oplus \mathbb{Z}_6$. We see that only G_1 contains elements of order 36 and is the only cyclic group. Hence G_1 cannot be isomorphic to G_2 , G_3 or G_4 . Of the three groups G_2 , G_3 and G_4 , only G_2 contains elements of order 18, so G_2 cannot be isomorphic to G_3 or G_4 . Finally, G_3 has elements of order 12 but G_4 does not, so G_3 and G_4 are not isomorphic, and we have now shown that no two of the four groups are isomorphic to each other.

As a slightly harder example, $\mathbb{Z}_2 \oplus \mathbb{Z}_2 \oplus \mathbb{Z}_4$ is not isomorphic to $\mathbb{Z}_4 \oplus \mathbb{Z}_4$, because the former has 7 elements of order 2, whereas the latter has only 3.