

More Linear Algebra¹

Chris Godsil
Combinatorics & Optimization
University of Waterloo

©2019

¹version: January 29, 2021

Preface

These notes are meant to provide an introduction to fundamental parts of linear algebra, as might be applied to problems in combinatorics. I assume the reader has had a first course in linear algebra, and is familiar with determinants.

To Do

1. Interlacing, via Courant-Fischer and by rational functions. Need equitable partitions for this.
2. Walk modules. Controllable graphs.
3. Quadrics.
4. Matrix exponential.
5. Lie algebras, $sl(2)$ reps.
6. Perron-Frobenius, symbolic dynamics.
7. Perturbation theory.

Contents

1	Spaces and Subspaces	1
1.1	Vector Spaces	1
1.2	Subspaces	3
1.3	Linear Mappings	4
1.4	Duals and Adjoint	6
1.5	Bilinear Forms	8
1.6	Counting	9
1.7	Normal Forms	11
1.8	Groebner Bases	13
1.9	Codes	15
2	Primary Decomposition	19
2.1	Modules	19
2.2	Control Theory	20
2.3	Sums	22
2.4	Invariant Sums	23
2.5	Minimal Polynomials	25
2.6	Primary Decomposition	27
2.7	The Degree of the Minimal Polynomial	29
2.8	Root Spaces	30
2.9	Examples of Root Spaces	31
2.10	Differential Equations	33
2.11	Linear Recurrence Equations	35
2.12	Diagonalizability	36
3	Frobenius Normal Form	39
3.1	Companion Matrices	39
3.2	Transposes	41

3.3	Eigenvectors for Companion Matrices	42
3.4	Inverses of Companion Matrices	44
3.5	Cycles	46
3.6	Circulants and Cyclic Codes	47
3.7	Frobenius Normal Form	49
3.8	Applications	52
3.9	Nilpotent Matrices	53
3.10	A Similarity Condition	55
3.11	Triangular Maps	57
3.12	Triangulations	59
3.13	The “Fundamental” “Theorem of Algebra”	59
3.14	The Kronecker Product	63
4	Orthogonality	67
4.1	Properties of Projections	67
4.2	Matrices Representing Projections	68
4.3	Least Squares	70
4.4	Orthogonal Polynomials	71
4.5	The Three-Term Recurrence	72
4.6	Numerical Integration	73
5	Eigenthings	77
5.1	Self-Adjoint Operators	77
5.2	Diagonalizability	78
5.3	Diagonalizability, Again	79
5.4	Eigenvectors and Optimization	80
5.5	The Singular Value Decomposition	81
5.6	Least Squares	83
5.7	Legendre Polynomials	85
5.8	Computing Eigenvalues	87
5.9	Jacobi: An Example	89
6	Spectral Decomposition	91
6.1	Self-Adjoint Operators	91
6.2	Commutative Algebras	92
6.3	Normal Operators	94

7 Norms	97
7.1 Convexity	97
7.2 Extreme Points	99
7.3 Norms	100
7.4 Dual Norms	102
7.5 Matrix Norms	103
7.6 Examples	105
7.7 Matrix Functions	106
7.8 Powers	108
7.9 Contractions	110
7.10 Projections	112
7.11 Contractions	114
7.12 Perron	116
8 Geometry	121
8.1 Semilinear Forms	121
8.2 The Classification of Forms	122
8.3 Gram Matrices	123
8.4 Equiangular Lines	124
8.5 Tight Frames	125
8.6 Another Gram Matrix	126
8.7 The Orthogonal Group	127
8.8 Skew-Symmetric to Orthogonal	128
8.9 Reflections	129
9 Positive Semidefinite Matrices	133
9.1 Factorizing Positive Semidefinite Matrices	133
9.2 Computing Cholesky	135
9.3 Polynomial Examples	136
9.4 Positive Semidefinite Matrices	138
10 Tensors	141
10.1 Tensor Products	141
10.2 Quadratic Tensors	142
10.3 Cubic Tensors	144
10.4 Multiplication	146
10.5 Semifields	147

11 Control	149
11.1 Buffalos	149
11.2 Burgers	151
11.3 Controllability	152
11.4 Observability	155
11.5 Feedback and Controllability	156
11.6 Canonical Forms	158
11.7 Eigenvalues and Controllability	159
11.8 Observers	161
11.9 Transfer Matrices	162
12 The Smith Normal Form	165
12.1 Domains	165
12.2 Localization	167
12.3 Fitting	168
12.4 Hermite	169
12.5 Smith Normal Form	171
13 Polynomial and Rational Matrices	175
13.1 Series	176
13.2 Polynomial Matrices	177
13.3 Paraunitary Matrices	180
13.4 Division	182
13.5 Cayley-Hamilton	183
13.6 Greatest Common Divisors	185
13.7 An Identity	187
13.8 Resolvents	189
14 Determinants	195
14.1 Permutations	195
14.2 The Sign of a Permutation	197
14.3 Permutation Matrices	198
14.4 Definition of the Determinant	199
14.5 The Determinant is Multiplicative	201
14.6 The Laplace Expansion	206
14.7 The Characteristic Polynomial of a Matrix	208
14.8 An Algorithm	210
14.9 Summary	213

14.10 Groups	214
15 Rings, Fields, Algebras	217
15.1 Rings	217
15.2 Fields	218
15.3 Algebras	219

Chapter 1

Spaces and Subspaces

We review the basic results on vector spaces.

1.1 Vector Spaces

We assume familiarity with the basic terminology of vector spaces—linear combinations, subspaces, linear dependence and independence, span, spanning sets, and bases. We present a proof of the existence of bases (in vector spaces with a finite spanning set).

We define a *circuit* in a vector space V to be a minimal dependent set. Thus if C is a circuit and x is any element of C then $C \setminus x$ is linearly independent. Hence $C \setminus x$ and C have the same span.

1.1.1 Lemma. *If the vector v lies in the span of a set S , then there is a circuit in $S \cup v$ that contains v .*

Proof. Suppose that v is a linear combination of the vectors x_1, \dots, x_k from S , and that v is not a linear combination of any subset of S with fewer than k elements. Then x_1, \dots, x_k is linearly independent, for otherwise it contains a circuit and by deleting an element of this circuit, we obtain a set of $k - 1$ vectors whose span still contains v . It follows that if for some i , the set

$$\{v, x_1, \dots, x_k\} \setminus x_i$$

is linearly dependent, then v is a linear combination of at most $k - 1$ elements of S . Therefore this set is linearly independent for each i , and so we conclude that $\{v, x_1, \dots, x_k\}$ is a circuit. \square

A basis, we recall, is a linearly independent spanning set. We show they exist if V has a finite spanning set S . If S is linearly independent, there is nothing to prove. Otherwise S contains a circuit C ; if $x \in C$ then $C \setminus x$ and C have the same span, and consequently $S \setminus x$ and S have the span. Therefore, by deleting a finite number of elements from S , we obtain a linearly independent set S_1 with the same span as S , and so S_1 is our basis.

Now we show that all finite bases have the same size. To do this we prove the following:

1.1.2 Lemma. *Let V be a vector space. If S is a finite linearly independent subset of V and T is a spanning set, then $|S| \leq |T|$.*

Proof. We prove the result by induction on $|S \setminus T|$. Set k equal to $|S \setminus T|$; if $k = 0$ the result is immediate, so suppose $k > 0$. Choose a vector u from $S \setminus T$.

Since T is a spanning set, u is a linear combination of elements of T , and therefore by the lemma above there is a circuit C in $T \cup u$ that contains u . Since S is linearly independent, C is not contained in S and therefore there is an element v in C that does not lie in S . Now v lies in the span of $C \setminus v$, and

$$C \setminus v \subseteq (T \setminus v) \cup u.$$

Therefore v belongs to the span of $(T \setminus v) \cup u$. Since this span contains $T \setminus v$, it contains T .

We conclude that $(T \setminus v) \cup u$ is a spanning set in V that meets S in $k + 1$ elements. \square

It follows from this that any two finite linearly independent spanning subsets of V have the same size, which we define to be the *dimension* of V . A vector space has finite dimension if and only if it has a finite basis. If V has dimension n then any independent set of size n is a basis, as is spanning set of size n . Each independent set is contained in a basis and, as we already knew, each spanning set contains a basis.

If $\alpha = (v_1, \dots, v_n)$ is an ordered basis for the vector space V and $w \in V$ then there are unique scalars a_1, \dots, a_n such that

$$w = \sum_{i=1}^n a_i v_i.$$

The *coordinate vector* $[w]_\alpha$ of w with respect to α is the $n \times 1$ matrix with entries a_1, \dots, a_n . The function that maps w to $[w]_\alpha$ is an injective map from V to \mathbb{F}^n .

We can also show that

$$[w + x]_\alpha = [w]_\alpha + [x]_\alpha, \quad [cw]_\alpha = c[w]_\alpha.$$

This shows that the coordinate map is an example of a linear mapping.

1.2 Subspaces

The intersection of any two subspaces (indeed, of any collection of subspaces) is a subspace. The union of two subspaces is rarely a subspace—in particular no vector space over an infinite field can be expressed as the union of a finite set of proper subspaces. There is a replacement for union though: the sum $U + V$ of two subspaces U and V . We define this by

$$U + V := \{u + v : u \in U, v \in V\}.$$

We see that $U + V$ is the span of $U \cup V$ and therefore it is a subspace and it is contained in any subspace that contains U and V . Consequently it is the intersection of all subspaces that contained U and V and it follows that the subspaces of a vector space, with the operations of intersection and sum, forms a lattice. If $U \cap V = \{0\}$, we say that $U + V$ is the *direct sum* of U and V .

Here we are concerned with the dimension of $U + V$. For this we need some preliminaries. Suppose U is a subspace of W . We say that a subspace V of W is a *complement* to U if $U \cap V = \{0\}$ and $U + V = W$. We construct examples as follows. Suppose S is a basis of W and (S_1, S_2) is a partition of S into two parts. Let U_i denote the span of S_i . Then $U_1 + U_2$ contains S , and hence it is equal to W . It is also not hard to show that $U_1 \cap U_2 = \{0\}$. Hence U_2 is a complement to U_1 (and vice versa).

1.2.1 Lemma. *Let W be a vector space with finite dimension. Then any subspace of W has a complement.*

Proof. Let U be a subspace of W and let S be a basis for U . Then there is a basis T for W that contains S , let V be the span of $T \setminus S$. □

1.2.2 Theorem. *If U and V are finite-dimensional subspaces of V , then*

$$\dim(U + V) = \dim(U) + \dim(V) - \dim(U \cap V).$$

Proof. We first establish a special case of the theorem: if U_1 and U_2 are subspaces and $U_1 \cap U_2 = \{0\}$, then

$$\dim(U_1 + U_2) = \dim(U_1) + \dim(U_2).$$

To derive this, we note that if S_i is an independent subset of U_i ($i = 1, 2$) and $U_1 \cap U_2 = \{0\}$ then $S_1 \cup S_2$ is linearly independent. Hence the union of a basis of U_1 and a basis of U_2 is a basis for $U_1 + U_2$.

Now we consider the general case. Let V_1 be a complement to $U \cap V$ in V . Then by what we have just proved,

$$\dim(V_1) = \dim(V) - \dim(U \cap V).$$

We show that V_1 is a complement to U in $U + V$. First

$$U + V_1 = U + ((U \cap V) + V_1) = U + V.$$

Second, $U \cap V_1 \subseteq U \cap V$ and $U \cap V_1 \subseteq V_1$, so

$$U \cap V_1 \subseteq (U \cap V) \cap V_1 = \{0\}.$$

Therefore V_1 is a complement to U in $U + V$ and consequently

$$\dim(V_1) = \dim(U + V) - \dim(U).$$

The two expressions for $\dim(V_1)$ imply the result. □

1.3 Linear Mappings

Let V and W be vector spaces over the same field. A function T with domain U and codomain V is a *linear mapping* from U to V if, for all vectors u_1 and u_2 in U ,

$$T(u_1 + u_2) = T(u_1) + T(u_2)$$

and if, for all scalars c and all vectors u in U ,

$$T(cu) = cT(u).$$

To specify a linear mapping, we must explicitly give its codomain. (This matters most when we consider adjoints.)

A bijective linear mapping is called an *isomorphism*. All this should be familiar. The image and kernel of a linear mapping T are subspaces. The dimension of $\text{im}(T)$ is its *rank* and the dimension of $\ker(T)$ is its *corank*. The following important relation between these parameters is sometimes called the “dimension theorem” for linear mappings.

1.3.1 Theorem. *If T is a linear mapping with domain V then*

$$\text{rk}(T) + \text{cork}(T) = \dim(V).$$

Proof. Choose a basis v_1, \dots, v_n for V such that v_1, \dots, v_k is a basis for $\ker(T)$. Let U be the span of v_{k+1}, \dots, v_n . If $u \in U$ and $Tu = 0$, then

$$u \in U \cap \ker(T) = \{0\}.$$

Hence the set $T(v_{k+1}), \dots, T(v_n)$ is linearly independent, and consequently it is a basis for $\text{im}(T)$. \square

This is perhaps the most useful formula in linear algebra. An important consequence is that, if T maps V to itself, then it is onto if and only if it is one-to-one.

The coordinate map with respect to a basis is an important example of a linear mapping.

If A is an $m \times n$ matrix over \mathbb{F} then the function that sends $x \in \mathbb{F}^n$ to Ax in \mathbb{F}^m is a linear mapping, often denoted T_A . This gives an even more important class of examples. Note that $\ker(T_A)$ is the null space of A and $\text{im}(T_A)$ is the column space of A , so the dimension theorem yields that

$$\text{rk}(A) + \text{cork}(A) = n.$$

As an application, we rederive the formula for the dimension of the sum of two subspaces. If U and V are vector spaces over the same field, their *external direct sum* is the vector space with vectors

$$\{(u, v) : u \in U, v \in V\},$$

where

$$(u_1, v_1) + (u_2, v_2) = (u_1 + u_2, v_1 + v_2)$$

and

$$c(u, v) = (cu, cv).$$

We denote this by $U \oplus V$, and claim that

$$\dim(U \oplus V) = \dim U + \dim V.$$

Now suppose that U and V are subspaces of W . Then we can define a linear map S from $U \oplus V$ to W by

$$S: (u, v) \mapsto u - v.$$

Note that S is a linear map from $U \oplus V$ to the subspace $U + V$ of W . It is easy to see that S is onto, and that its kernel consists of the vectors (x, x) , where $x \in U \cap V$. Hence

$$\dim(U + V) = \text{rk}(S) = \dim(U) + \dim(V) - \dim(U \cap V).$$

Define

$$(U, 0) := \{(u, 0) : u \in U\}$$

and define $(0, V)$ similarly. Then $(U, 0)$ and $(0, V)$ are subspaces of $U \oplus V$ having zero intersection and

$$U \oplus V = (U, 0) + (0, V).$$

Thus an external direct sum is a direct sum of subspaces, as in the previous section.

The term “external direct sum” is somewhat confusing. It may help to view this as follows. We have a simple construction of a vector space W from two vector spaces U and V over a field \mathbb{F} . The space W is the direct sum, in our original sense, of subspaces isomorphic to U and V .

1.4 Duals and Adjoints

Since we can add linear transformations from V to W and multiply them by scalars, the set $\mathcal{L}(V, W)$ of all linear transformations from V to W forms a vector space. Hence:

1.4.1 Theorem. *If V and W are vector spaces over \mathbb{F} , then $\mathcal{L}(V, W)$ is a vector space with dimension $\dim(V) \dim(W)$.*

Proof. We present you a set of linear mappings, and invite you to prove they form a basis.

Let v_1, \dots, v_n be a basis for V and w_1, \dots, w_m be a basis for W . Let $E_{i,j}$ be the element of $\mathcal{L}(V, W)$ given by

$$E_{i,j}(v_r) = \begin{cases} w_j, & \text{if } r = i; \\ 0, & \text{otherwise.} \end{cases}$$

(We use the fact that a linear transformation can be defined by specifying its values on a basis.) This set of $\dim(V) \dim(W)$ operators is the subset we promised. \square

Here we will be most interested in the *dual space* $\mathcal{L}(V, \mathbb{F})$, which we denote by V^* . We consider some examples.

Suppose V is the space of all polynomials over \mathbb{F} . If $\psi \in V^*$, then ψ is determined by its values on a basis, and hence determined by its values on the powers of x . If we denote $\psi(x^n)$ by ψ_n , then we find that

$$\psi : \sum_{i=0}^m p_i x^i \rightarrow \sum_{i=0}^m p_i \psi_i.$$

Thus each sequence $(\psi_n)_{n \geq 0}$ determines an element of V^* . It follows that we can identify V^* with the space of all formal power series in x .

Each element v of V gives rise to a map from V^* to \mathbb{F} , that sends ψ in V^* to $\psi(v)$ in \mathbb{F} . This map is linear and injective, and allows us to identify V with a subspace of $(V^*)^*$. The previous example shows that this map need not be an isomorphism in general, but it is an isomorphism when $\dim(V)$ is finite. (This follows from the observation that V , V^* and V^{**} all have the same dimension.)

If $V = \mathbb{F}^n$, then the map that sends an element v to its i -th coordinate is linear, and so belongs to V^* . In this case $V^* \cong V$.

If $V = \text{Mat}_{n \times n}(\mathbb{F})$, then the trace function is an element of V^* .

We cannot resist remarking on one special property of V^* . There is a natural product on it: if $f, g \in V^*$ then fg is defined by $(fg)(u) = f(u)g(u)$.

Let T be a linear map from V to W . If $g \in W^*$, then the composition $g \circ T$ is a linear mapping from V to \mathbb{F} ; hence it is an element of V^* . Thus we have a mapping that takes an element g of W^* to an element $g \circ T$ in V^* . This map is linear (prove it!), and is called the *adjoint* of T . We denote it by T^* . (We also offer a warning: if T is a linear operator on an inner product space, the term 'adjoint' is applied to a different map.)

(1) Prove that T^* is linear.

- (2) Prove that T is one-to-one if and only if T^* is onto, and that T is onto if and only if T^* is one-to-one.
- (3) Prove that $T^{**} = T$.
- (4) Prove that V is isomorphic to a subspace of V^{**} .

1.5 Bilinear Forms

Suppose Φ is a linear mapping from V to V^* . If $u, v \in V$, then the map

$$(u, v) \mapsto \Phi(u)(v)$$

is linear in each variable. Such a map is called a *bilinear form*. The simplest example arises if we take V to be the space of $n \times 1$ matrices over \mathbb{F} . Then we can identify V^* with the space of $1 \times n$ matrices. If $v^T \in V^*$ and $u \in V$, then the value of v^T on u is $v^T u$. So we may take Φ to be the transpose map, and then the bilinear form takes (u, v) to $u^T v$. We generally denote the value of a bilinear form by $\langle u, v \rangle$.

If $u \in V$ and $\Phi(u)(v) = 0$ for all v then $\Phi(u)$ must be the zero vector, and so $u \in \ker \Phi$. If $\Phi(u)(v) = 0$ for all u , then $\text{im } \Phi$ lies in the subspace of V^* formed by the elements f such that $f(v) = 0$. If V is finite dimensional, then V and V^* have the same dimension and $\ker \Phi$ is the zero subspace if and only if $\text{im } \Phi = V^*$. We say that a bilinear form is *non-degenerate* if Φ is invertible; in this case Φ is an isomorphism and we have the following description of V^* :

1.5.1 Lemma. *Let V be a finite-dimensional vector space with a non-degenerate bilinear form. If $f \in V^*$, then there is a vector v in V such that $f(x) = \langle v, x \rangle$. \square*

A bilinear form is symmetric if

$$\langle u, v \rangle = \langle v, u \rangle$$

for all u and v . It is *alternating* if

$$\langle u, v \rangle = -\langle v, u \rangle$$

and $\langle u, u \rangle = 0$ for all u . (The first condition implies the second unless we are working over a field of characteristic two.)

We describe one simple construction of bilinear forms. Let A be an $n \times n$ matrix over \mathbb{F} . If u and v belong to \mathbb{F}^n , define

$$\langle u, v \rangle := u^T A v.$$

It is easy to verify this is bilinear. It is non-degenerate if and only if A is invertible. It is symmetric if and only if $A = A^T$ and alternating if and only if both $A^T = -A$ and all diagonal entries of A are zero.

If S is a subset of V then we define S^\perp to be the set of vectors v such that $\langle v, x \rangle = 0$ for all x in S . (In practice, S will usually be a subspace or a vector.) It is true that if U is a subspace of V , then

$$\dim U^\perp = \dim V - \dim U;$$

but we leave you to prove this. (See the exercises at the end of this section.)

(1) If U is a subspace of V , show that $V = U + U^\perp$ if and only if $U \cap U^\perp = \{0\}$.

(2) Given that $\dim(U^\perp) = \dim(V) - \dim(U)$, prove that $U^{\perp\perp} = U$.

1.6 Counting

We count bases and subspaces in vector spaces over $GF(q)$. Throughout this section we assume that \mathbb{F} has order q . Let $V = \mathbb{F}^n$. Then V contains exactly q^n elements.

We begin by counting the number of subspaces of dimension 1. First we note that two distinct subspaces of dimension 1 have only the zero vector in common, and that a subspace of dimension 1 contains exactly $q - 1$ non-zero vectors. It follows that there are exactly $(q^n - 1)/(q - 1)$ 1-dimensional subspaces of V . This number plays quite a role in what follows, so we define

$$[n] := \frac{q^n - 1}{q - 1}.$$

(We will write $[n]_q$ if we need to make the order of \mathbb{F} explicit.) Note that $[1] = 1$ and $[2] = q + 1$.

We next determine the number of ordered k -tuples (v_1, \dots, v_k) of vectors from V such that v_1, \dots, v_k is linearly independent. Suppose we have such a $(k - 1)$ -tuple. We can extend it to a k -tuple by choosing vector not in the $(k - 1)$ -dimensional subspace spanned by the $(k - 1)$ -tuple. There are $q^n - q^{k-1}$ such

factors, and now a simple induction argument yields that the number of ordered k -tuples of linearly independent vectors is

$$(q^n - 1) \cdots (q^n - q^{k-1}) = q^{\binom{k}{2}} (q - 1)^k [n][n-1] \cdots [n-k+1].$$

Since each k -tuple of linearly independent vectors spans a unique subspace of dimension k , and since each subspace of dimension k gives rise to exactly

$$q^{\binom{k}{2}} (q - 1)^k k[k-1] \cdots [1]$$

k -tuples of linearly independent vectors, we find that the number of subspaces of dimension k is

$$q^{\binom{k}{2}} (q - 1)^k [n][n-1] \cdots \frac{[n-k+1]}{q^{\binom{k}{2}}} (q - 1)^k k[k-1] \cdots [1] = \frac{[n][n-1] \cdots [n-k+1]}{[k][k-1] \cdots [1]}.$$

This suggests the use of the following notation. We define

$$[n]! := [n][n-1] \cdots [1]$$

and

$$\begin{bmatrix} n \\ k \end{bmatrix} := \frac{[n]!}{[k]![n-k]!}. \quad (1.6.1)$$

The right side of (1.6.1) is known as the *Gaussian binomial coefficient*. Using it, we have:

1.6.1 Theorem. *The number of subspaces of dimension k in a vector space of dimension n over a field of order q is $\begin{bmatrix} n \\ k \end{bmatrix}$.* \square

We note another consequence. An ordered basis for \mathbb{F}^n is the same thing as an invertible $n \times n$ matrix. Hence:

1.6.2 Lemma. *The number of invertible $n \times n$ matrices over a field of order q is $q^{\binom{n}{2}} (q - 1)^n [n]!$.* \square

Although it may not be immediately apparent, the Gaussian binomial coefficient is a polynomial in q .

- (1) Prove the recurrence for Gauss...
- (2) Let U be a fixed subspace of \mathbb{F}^n with dimension k . Compute the number of ℓ -dimensional subspaces V of \mathbb{F}^n such that $V \cap U = \{0\}$.
- (3) Let U and V be subspaces of \mathbb{F}^n such that $\dim(U) = k$, $\dim(V) = n - k$ and $U \cap V = \{0\}$, where $2k \leq n$. Compute the number of subspaces W with dimension k such that

$$W \cap U = W \cap V = \{0\}.$$

1.7 Normal Forms

There are a number of cases where we wish to decide if two linear independent sets in a vector space span the same subspace. We can answer this by constructing normal forms.

The conventional way to do this is using reduced row echelon form. At first glance this only solves our problem for the vector spaces \mathbb{F}^n , but we can always translate our problem to \mathbb{F}^n using the coordinate map with respect to some basis.

Let V be a finite dimensional vector space and let f_1, \dots, f_n be an ordered basis for V^* . If $v \in V$, we define the *height* of v to be the least value of i such that $f_i(v) \neq 0$. We write $v < w$ if the height of v is less than the height of w . (Technically the relation we have just introduced is a pre-order on V .)

1.7.1 Lemma. *Let S be a finite subset of the finite-dimensional vector space V . Then there is a subset T of V such that $\text{span}(S) = \text{span}(T)$ and no two elements of T have the same height.*

Proof. We convert S to T in a number of steps, as follows. If S does not contain a pair of elements of the same height, there is nothing to be done. If S does contain a pair of elements of the same height, choose a pair v and w with the least possible height, i say. Then replace w by

$$w' := w - \frac{f_i(v)}{f_i(w)}v.$$

Then the height of w' is greater than the height of w and

$$\text{span}((S \setminus w) \cup w') = \text{span}(S).$$

In this we may replace each element of height i , other than v , by a vector with greater height, eventually reaching the situation where v is the unique element of height i .

By induction on the number of elements, we may convert $S \setminus v$ to a set T_1 such that $\text{span}(S \setminus v) = \text{span}(T_1)$, no two elements of T_1 have the same height, and the least height of an element of T_1 is greater than i . Then we may take T to be $T_1 \cup v$. \square

We remark that this lemma holds in the vector space of all polynomials. Take f_r to be the element of the dual space that maps a polynomial to the coefficient of x^r . Then $f_r(p) = 0$ for all r if and only if p is the zero polynomial, and

this is all we need to make the argument work. One interesting point is that the f_r 's do not form a basis of the dual space. Note also that in this case the height of a polynomial is just its degree. For $\text{Pol}(\mathbb{R})$, the above lemma asserts that each subspace with finite dimension has a spanning set containing at most one polynomial of each degree.

We will say that S is in *echelon form* if it contains at most one element of each height. If S is in echelon form, then $|S| \leq n$ and S is linearly independent. Further, if S does not contain an element of height i , then no element of $\text{span}(S)$ has height i . Consequently a basis for V must have size n .

We will say S is in *normal form* if:

- (a) If $v \in S$ has height i , then $f_i(v) = 1$.
- (b) There is at most one element of any given height in S .
- (c) If v and w are elements of S with heights i and j respectively and $i < j$, then $f_j(v) = 0$.

These properties are a straightforward extension of the concept of reduced row-echelon form.

1.7.2 Lemma. *Let V be a finite dimensional vector space, and suppose S_1 and S_2 are subsets of V in normal form. If S_1 and S_2 have the same span, they are equal.*

Proof. Suppose $U = \text{span}(S_1)$ and let H denote the set of heights of elements of S_1 . If $i \in H$, let v_i denote the element of S_1 with height i . If $x \in V$, let z be given by

$$z := x - \sum_{i \in H} f_i(x) v_i.$$

Then $f_j(z) = 0$ when $j \in H$. If $x \neq z$, then the height of $x - z$ is not an element of H . Therefore $x \in U$ if and only if $x = z$.

Note that i is the height of an element in $\text{span}(S_1)$ if and only if $i \in H$. Accordingly H is also the set of heights of elements of S_2 . Let w_1 and w_2 denote the elements of least height in S_1 and S_2 respectively. Then w_1 and w_2 have the same height, k say.

Since $w_2 \in U$, so is $w_1 - w_2$. But

$$f_i(w_1 - w_2) = 0$$

for all i in H , and therefore $w_1 = w_2$.

The set of vectors in U with height greater than k is a subspace of U ; the heights of the elements in this subspace are elements of H that are greater than k . As both $S_1 \setminus w_1$ and $S_2 \setminus w_2$ are bases for this subspace in normal form, it follows by induction (on the dimension of U) that they are equal.

1.8 Groebner Bases

Let R be the ring of polynomials $\mathbb{F}[x_1, \dots, x_n]$. An *ideal* in R is a subset I of R such that, if $p \in R$ and $q \in I$, then $pq \in I$. An ideal is a subspace of R that is invariant under multiplication by elements of R . An ideal is *finitely generated* if it contains polynomials g_1, \dots, g_r such that

$$I = \left\{ \sum_{i=1}^r p_i g_i, p_i \in R \right\}.$$

If $\alpha = (a_1, \dots, a_n)$ is a sequence of non-negative integers, then

$$x^\alpha := \prod_{i=1}^n x_i^{a_i}.$$

We use $|\alpha|$ to denote $\sum_i a_i$, and we say that x^α is a *monomial* of degree $|\alpha|$. A *monomial order* is a total order ' \leq ' of the set of all monomials in x_1, \dots, x_n such that

- (a) For any λ , we have $1 \leq x^\lambda$.
- (b) If α, β and γ are non-negative and $x^\alpha \leq x^\beta$, then $x^{\alpha+\gamma} \leq x^{\beta+\gamma}$.

One important property of monomial orders is that any set of monomials has a least element. It follows from (b) that a monomial order is an extension of the usual partial order on monomials, that is, if x^α divides x^β , then $x^\alpha \leq x^\beta$. We define the *height* of a monomial to be the number of monomials less than it, relative to the given order. The height of a polynomial is the height of its leading term. (Note that this is consistent with our usage in the previous section—the map that takes a polynomial to the coefficient of a given monomial in it lies in the dual space.)

Once we have chosen a monomial order, each polynomial in R has a unique *leading term* of the form cx^α , for some monomial x^α . The scalar c is the leading coefficient of the polynomial and x^α is its leading monomial. If I is an ideal of R , the *leading-term ideal* of I is the ideal generated by the leading terms of the polynomials in I .

1.8.1 Lemma. *Any ideal of $\mathbb{F}[x_1, \dots, x_n]$ that is generated by monomials is generated by a finite set of monomials.*

Proof. Let I be an ideal of R , let \mathcal{M} denote the set of leading monomials of I , and let J be the ideal generated by \mathcal{M} . If \mathcal{M} contains both x^α and x^β and x^α divides x^β , then $\mathcal{M} \setminus x^\beta$ generates J . Hence we may assume that no element of \mathcal{M} divides another and so, by Lemma 1.8.3, it follows that \mathcal{M} is finite. \square

1.8.2 Theorem. *Every ideal of $\mathbb{F}[x_1, \dots, x_n]$ is finitely generated.*

Proof. Let I be an ideal of R , let \mathcal{M} denote the set of leading monomials of I , and let J be the ideal generated by \mathcal{M} . Let g_1, \dots, g_m be set of polynomials in I whose leading terms generate the leading-term ideal of I . We show that these polynomials generate I .

If $h \in I$ then there are polynomials a_1, \dots, a_m such that the leading term of

$$\sum_{i=1}^m a_i g_i$$

equals the leading term of h . Hence $h - \sum_{i=1}^m a_i g_i$ is a polynomial in I with leading term less than the leading term of h . By induction on the height, we conclude that this polynomial lies in the ideal generated by the g_i 's, and hence that $h \in I$. \square

A *Groebner basis* for an ideal in $\mathbb{F}[x_1, \dots, x_n]$ is a generating set g_1, \dots, g_n such that the leading terms of the g_i 's generate the leading-term ideal of I . The previous result is a nice application of Groebner bases, which are a very useful tool in computational algebra. We have managed to sidestep the question of how we might actually find a Groebner base of an ideal. (They are not unique.)

If x and y are vectors over an ordered ring, e.g., the integers, we write $x \geq y$ if the entries of $x - y$ are all non-negative.

1.8.3 Lemma. *If S is an infinite set of n -tuples of non-negative integers, it contains a pair of distinct elements x and y such that $x \leq y$.*

Proof. We proceed by induction on n ; the case $n = 1$ is trivial. Consider the subset $S(i, j)$ of S defined by

$$S(i, j) := \{x \in S : x_i = j\}.$$

If $S(i, j)$ is infinite for some i and j , then the set of $(n - 1)$ -tuples we get by dropping the i -th coordinate of each element of $S(i, j)$ contains a comparable pair of distinct elements.

So we may assume that $S(i, j)$ is finite for all i and j . Choose an element z of S . Then

$$\bigcup_{i=1}^n \bigcup_{j \leq z_i} S(i, j)$$

is a finite union of finite sets. Accordingly there are infinitely many elements y in S such that $y \geq z$. \square

1.8.4 Corollary. *Let A and B be integer matrices and let \mathcal{N} denote the set of non-negative integer solutions to $AX = B$. Define an element of \mathcal{N} to be indecomposable if it cannot be written as the sum of two non-zero elements of \mathcal{N} . Then the number of indecomposable elements of \mathcal{N} is finite.* \square

1.9 Codes

The goal of coding theory is to successfully transmit information, despite the occurrence of errors in transmission. We describe one way of doing this.

We assume we have a message to transmit, which is a long binary string. We divide the message into *input words* of length k , padding the message if needed. If x is an input word, we transmit

$$w^T = x^T G. \tag{1.9.1}$$

If the rows of G are linearly dependent, there will be vectors x such that $x^T G = 0$ but $x \neq 0$. This is clearly no use to us, so we assume that $\text{rk}(G) = k$.

In place of w the receiver actually receives $(w + e)^T$, where e is an error vector. The receiver's task to determine e , and then to solve the equation (1.9.1) for x . The possible form of e depends very much on the details. We will assume that each bit is sent correctly with probability $1 - p$, and with probability p it is changed from a 0 to a 1, or vice versa. The expected number of errors when we send n bits is np . Generally p is small.

An input word is just an element of \mathbb{Z}_2^k . The words we transmit, which we call *code words*, are elements of $\text{row}(G)$. If the error vector $e \notin \text{row}(G)$ then $(w + e)^T \notin \text{row}(G)$. We can decide if a received word is a code word by attempting to solve (1.9.1) for x . Define the *Hamming distance* $h(v, w)$ between two

binary vectors to be the number of positions at which they differ. The *minimum distance* of a code is the minimum distance between two distinct code words. Our basic problem is to choose G so that its rank is large, and the minimum distance is at least $np + \gamma\sqrt{np}$, for some positive constant γ . (For example, $\gamma = 10$.)

We consider an example. Let G be given by

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

Thus we have 16 code words in our code C . The receiver needs a method to decide if a given element w of \mathbb{Z}_2^n lies in C . A convenient approach is to use a so-called *parity check matrix*. This is a $(n - k) \times n$ matrix H with linearly independent rows, such that

$$HG^T = 0.$$

Thus the null space of H is equal to the column space of G^T and a vector w^t lies in $\text{row}(G)$ if and only if $Hw = 0$. If we receive w and $Hw = 0$ then the first four entries of w are the input word. (Thus decoding is trivial.) Here we may take H as follows (see Exercise ?? below):

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

Suppose we send $w^T = x^T G$ and the $w + e_i$ arrives at the receiver. As usual, e_i is the i -th standard basis vector—we are assuming that exactly one error occurs in transmission. Then the receiver computes

$$Hw = H(w + e_i) = Hw + He_i = He_i.$$

This is the i -th column of H . Inspection reveals that the seven columns of H are distinct, and thus we can decide which entry of w has been corrupted in transmission. Consequently we can correct any single error.

We note that since the columns of H are distinct, if $i \neq j$, then

$$H(e_i + e_j) = He_i + He_j \neq 0;$$

therefore we can detect if one or two transmission errors occur.

Thus we might choose simply to detect errors, and request a retransmission whenever an error occurred. Or we might know that the probability of two errors in one word is too low to worry about, and correct each error that we noted. In this case, if two errors did occur, our “correction” would not give us the word that was really sent.

Note that we cannot hope to correct or even detect all errors, but we can aim to reduce the proportion of undetected or wrongly corrected errors to a specified level.

(1) If

$$G := (I \quad G_1)$$

and

$$H := (G_1^T \quad -I),$$

show that the null space of H equals $\text{row}(G^T)$.

(2) Suppose G is a $k \times n$ matrix in reduced row echelon form and let G_1 be the $k \times (n - k)$ matrix we get from G by deleting the basic columns of G . Construct an $(n - k) \times n$ matrix as follows: if the j -th column of G is the i -th non-basic column of G , the j -th column of H is $-e_i$; if the j -th column of G is basic and $G_{i,j} = 1$, then the j -th column of H is times the i -th column of G_1^T . Prove that $HG^T = 0$.

(3) Give a generator matrix G of order 11×15 such that the associated code can be used to correct any single error.

Chapter 2

Primary Decomposition

We use the primary decomposition to decompose vector spaces and linear mappings.

2.1 Modules

Let V be a vector space over \mathbb{F} and let T be an endomorphism of V . A subspace U of V is T -invariant if $Tu \in U$ for all elements u of U . If U is T -invariant, it is invariant under all matrices in the ring $\mathbb{F}[T]$ of polynomials in T . Hence it is a module over this ring; we may also refer to it as a T -module.

- (1) If $T = I$ then a T -invariant subspace is just another name for a subspace.
- (2) The zero subspace and V itself are T -invariant, for any T .
- (3) $\ker(T)$ is T -invariant, because if $u \in \ker(T)$ then $Tu = 0$, and certainly $0 \in \ker(T)$.
- (4) The range of T is T -invariant. For if u lies in the range of T then Tu is contained in the range of T .
- (5) If U is a subspace of V , the *preimage* of U relative to T is the set

$$\{v \in V : Tv \in U\}.$$

If U is T -invariant, then so is its preimage relative to T . (Since $\ker(T)$ is the preimage of $\{0\}$, this shows that $\ker(T)$ is T -invariant.)

(6) The intersection and sum of T -invariant subspaces are T -invariant.

If U is a T -invariant subspace, then $T|U$ denotes the endomorphism of U that is defined by

$$(T|U)(u) = Tu,$$

for all u in U . We call $T|U$ the *restriction* of T to U . If U is a 1-dimensional T -invariant subspace and u spans U , then Tu must be a scalar multiple of u . If u is a non-zero vector and $Tu = \theta u$, we say that u is an *eigenvector* of T with *eigenvalue* θ .

If $u \in W$, then the subspace spanned by vectors

$$T^r v, \quad r = 0, 1, \dots$$

is easily seen to be T -invariant. We call it the T -invariant subspace *generated* by v , and observe that is the smallest T -invariant subspace of W that contains v . A T -invariant subspace generated by a single vector u is called a *cyclic subspace* for T . Cyclic subspaces are perhaps the most important class of invariant subspaces.

(1) If $T \in \text{End}(V)$ and T is invertible, show that a T -invariant subspace is T^{-1} -invariant.

2.2 Control Theory

Consider a system of $n + 1$ bodies arranged in a line. Assume that if the temperature of the i -th body ($1 \leq i \leq n$) at time r is $t_i(r)$, then its temperature at time $i + 1$ is given by

$$t_i(r + 1) = \frac{1}{4}(t_{i-1}(r) + 2t_i(r) + t_{i+1}(r))$$

The temperature of the 0-th body is entirely under our control, we denote its value at time r by $u(r)$. The temperature of the $(n + 1)$ -st is fixed at zero. If $t(r)$ is the vector in \mathbb{R}^n with i -th entry $t_i(r)$ then t is determined by the equation of the form:

$$t(r + 1) = At(r) + u(r)b$$

and the temperature vector $t(0)$ at time zero. In particular, if $n = 4$ then

$$A = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0.25 & 0.5 & 0.25 & 0 & 0 & 0 \\ 0 & 0.25 & 0.5 & 0.25 & 0 & 0 \\ 0 & 0 & 0.25 & 0.5 & 0.25 & 0 \\ 0 & 0 & 0 & 0.25 & 0.5 & 0.25 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}.$$

By choosing different values for the terms of the sequence

$$u(0), u(1) \dots u(m)$$

we can reach a variety of different temperature distributions; are there any we cannot reach?

To study this we assume that $t(0) = 0$. Then

$$t(1) = u(0)b$$

$$t(2) = u(1)b + u(0)Ab$$

$$t(3) = u(2)b + u(1)Ab + u(0)A^2b$$

and

$$t(r+1) = \sum_{i=0}^r u(r-i)A^i b.$$

If W_r is the matrix

$$W_r = (b \quad Ab \quad \dots \quad A^{r-1}b)$$

then we see that

$$t(r+1) = W_r \begin{pmatrix} u(0) \\ \vdots \\ u(r) \end{pmatrix}.$$

The state $t(r+1)$ is therefore reachable if and only if it lies in the column space of W_r . When $r \geq n$, this column space is precisely the A -cyclic subspace generated by b . (As the vectors $A^r b$ lie in \mathbb{R}^n we have that A^n lies in the column space of W_n and, in general, the rank of W_m equal to the rank of W_n , whenever $m \geq n$.)

In our particular example above, W_6 is an upper triangular matrix with diagonal entries 4^{1-r} , for $r = 1, \dots, 6$. Therefore the cyclic subspace generated by b

is \mathbb{R}^6 , and so all states are reachable after at most six steps. If we change b to

$$\begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}$$

then the space of reachable states has dimension two—in this case all reachable states have $t_1(r) = t_6(r)$, $t_2(r) = t_5(r)$ and $t_3(r) = t_4(r)$.

- (1) Show that, even if $t(0) \neq 0$, it is still true that all states are reachable after at most n steps if and only if the A -cyclic subspace generated by b is all of \mathbb{R}^n .
- (2) A system given by

$$t(r+1) = At(r), \quad z(r) = c^T t(r)$$

is *observable* if, given $z(0), z(1), \dots, z(m)$ (where $m \geq n$) we can compute $t(i)$ for $i = 0, \dots, m$. (Note: the input for the computation is A , c and the values of z .) Show that this system is observable if and only if the cyclic A^T -subspace generated by c is equal to \mathbb{R}^n . Show further that, if the system is observable, we need at most n consecutive values of z to determine all previous states of the system.

2.3 Sums

We consider direct sums of subspaces. Suppose U_1, \dots, U_k are subspaces of V , and define U'_i to be the sum of the subspaces U_j , where $j \neq i$. We say that V is the *direct sum* of the subspaces U_i if V is the sum of the subspaces U_i and

$$U_i \cap U'_i = \{0\}, \quad (i = 1, \dots, k). \quad (2.3.1)$$

If this condition holds, we write

$$V = U_1 \oplus \dots \oplus U_k. \quad (2.3.2)$$

There is a condition equivalent to (2.3.1) that is often easier to work with: V is the direct sum of U_1, \dots, U_k if and only if for $i = 1, \dots, k-1$,

$$U_i \cap (U_{i+1} + \dots + U_k) = \{0\}.$$

We leave you to verify that these two conditions are equivalent.

As an easy consequence of the definition of direct sum, we have

$$\dim(V) = \dim(U_1) + \cdots + \dim(U_k).$$

There is a converse to this: if U_1, \dots, U_k are subspaces of V whose sum is V and

$$\sum_i \dim(U_i) = \dim(V),$$

then V is the direct sum of the U_i 's.

If (2.3.2) holds and $v \in V$, then v can be written in exactly one way as a sum

$$v = u_1 + \cdots + u_k,$$

where $u_i \in U_i$. Define a map $E_i: V \rightarrow U_i$ by $E_i(v) = u_i$. Then E_i is linear,

$$E_1 + \cdots + E_k = I,$$

and

$$E_i E_j = \begin{cases} E_i, & \text{if } i = j; \\ 0, & \text{otherwise.} \end{cases}$$

Note that the last condition implies that E_i is *idempotent*, that is, $E_i^2 = E_i$. We call E_i the *projection* onto U_i . Conversely, if E_1, \dots, E_k is a set of idempotents satisfying these conditions and U_i is the range of E_i , then V is the direct sum of the spaces U_i .

1. If u_1, \dots, u_n are elements of V and $U_i = \text{span}(u_i)$, show that V is the direct sum of U_1, \dots, U_n if and only if u_1, \dots, u_n is a basis for V .

2.4 Invariant Sums

If T is an endomorphism of V , we say a direct sum decomposition of V is *T-invariant* if each summand is. If V is the T -invariant direct sum of U_1, \dots, U_k and $v \in V$ then

$$v = u_1 + \cdots + u_k,$$

where $u_i \in U_i$. Hence

$$T(u) = (T \upharpoonright U_1)(u_1) + \cdots + (T \upharpoonright U_k)(u_k),$$

and so we say that T is the *direct sum* of the operators $T|_{U_i}$. It can be extremely useful to be able to decompose V into a T -invariant direct sum.

We develop a characterization of invariant direct sums in terms of projections. We use the following simple tool.

2.4.1 Lemma. *If E is idempotent, then $x \in \text{im}(E)$ if and only if $x = Ex$.*

Proof. If $x \in \text{im}(E)$ then $x = Ey$ for some y and therefore

$$Ex = E^2y = Ey = x.$$

If $x = Ex$ then clearly $x \in \text{im}(E)$. □

2.4.2 Theorem. *Suppose $V = V_1 \oplus \cdots \oplus V_k$ and let E_1, \dots, E_k be the set of projections corresponding to the subspaces V_i . Let T be a linear operator on V . Then the direct sum decomposition of V is T -invariant if and only if $TE_i = E_iT$ for each i .*

Proof. We first claim that if E is an idempotent then $\text{im}(E)$ is T -invariant if and only if $(I - E)TE = 0$.

Now $(I - E)TE = 0$ if and only if T maps $\text{im}(E)$ into $\ker(I - E)$. But $(I - E)x = 0$ if and only if $x = Ex$ and so the previous lemma implies that $\ker(I - E) = \text{im}(E)$. This proves our claim.

It follows from this claim that $\text{im}(I - E)$ is T -invariant if and only if $ET(I - E) = 0$.

If $TE = ET$, then both $ET(I - E)$ and $(I - E)TE$ are zero. Conversely, if

$$ET(I - E) = (I - E)TE = 0$$

then

$$0 = ET(I - E) - (I - E)TE = ET - TE$$

and so T and E commute. Hence we have shown that $\text{im}(E)$ and $\text{im}(I - E)$ are T -invariant if and only if $ET = TE$.

Let V'_i be the sum of the subspaces V_j for $j \neq i$. Then $V'_i = \text{im}(I - E_i)$, and so V_i and V'_i are both T -invariant if and only if E_i commutes with T . The theorem follows directly from this. □

Our next result identifies one case where we can express V as a sum of T -invariant subspaces.

2.4.3 Lemma. *Let T be an endomorphism of V . Then $V = \text{im}(T) + \text{ker}(T)$ if and only if $\text{im}(T) \cap \text{ker}(T) = \{0\}$.*

Proof. Suppose $\text{rk}(T) = k$ and $\dim(W) = n$. Then $\dim(\text{ker}(T)) = n - k$ and so $\text{im}(T) + \text{ker}(T) = n$ if and only if $\text{im}(T) \cap \text{ker}(T) = \{0\}$. \square

The constraint on T here may also be expressed thus: if $w \in W$ and $T^2 w = 0$ then $T w = 0$.

As an application of this lemma, suppose that T is idempotent. If $T^2 v = 0$, then $T v = 0$ and so no non-zero vector $T v$ lies in $\text{ker}(T)$. Hence V is the direct sum of $\text{im}(T)$ and $\text{ker}(T)$. Note that $T|_{\text{ker}(T)}$ is the zero map.

1. Show that a square matrix of the form

$$P := \begin{pmatrix} 0 & X \\ 0 & I \end{pmatrix}$$

is idempotent. If T is represented by the matrix

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix},$$

show that T fixes $\text{ker } P$ if and only if $C = 0$ and that T fixes $\text{col}(P)$ if and only if

$$XCX - AX + XD - B = 0.$$

2.5 Minimal Polynomials

Let T be an endomorphism of the n -dimensional vector space V . If $v \in V$, then there is a least positive integer r such that $T^r v$ lies in the span of $v, T v, \dots, T^{r-1} v$. Hence there are scalars a_1, \dots, a_r such that

$$T^r v + a_1 T^{r-1} v + \dots + a_0 v = 0.$$

It follows that there is a monic polynomial $\varphi(t)$ such that $\varphi(T)v = 0$. If φ_1 and φ_2 are two polynomials such that $\varphi_i(T)v = 0$, then for all polynomials $a_1(t)$ and $a_2(t)$,

$$(a_1(T)\varphi_1(T) + a_2(T)\varphi_2(T))v = 0,$$

from which it follows that if $\varphi(t)$ is the gcd of $\varphi_1(t)$ and $\varphi_2(t)$, then $\varphi(T)v = 0$. Consequently:

2.5.1 Lemma. Suppose T is an endomorphism of the finite-dimensional vector space V and $v \in V$. There is a unique monic polynomial of least degree $\psi_v(t)$ such that $\psi_v(T)v = 0$. The degree of ψ_v is equal to the dimension of the subspace generated by v . \square

We call $\psi_v(t)$ the *minimal polynomial* of T relative to v . Since $\dim V = n$, the degree of $\psi_v(t)$ is at most n .

Next we observe that space of endomorphisms of V has dimension n^2 , and therefore there is a least integer r , at most n^2 , such that I, T, \dots, T^r are linearly dependent. It follows that there is a unique monic polynomial ψ of least degree such that $\psi(T) = 0$. It is called the *minimal polynomial* of T . (If L_T denotes the linear operator on $\text{End}(V)$ given by $L_T(M) = TM$, then the minimal polynomial of T is the minimal polynomial of L_T relative to T itself.)

If $v \in V$, then certainly $\psi(T)v = 0$, and it follows that $\psi_v(t)$ must divide $\psi(t)$. Hence $\psi(t)$ is the least common multiple of the polynomials $\psi_v(t)$, as v runs over a basis of V .

2.5.2 Lemma. Suppose T is an endomorphism of the finite-dimensional vector space V and ψ is the minimal polynomial of T . Then each zero of ψ is an eigenvalue of T .

Proof. Suppose $\psi_v(\theta) = 0$. Then

$$\psi(t) = (t - \theta)\varphi(t)$$

and therefore

$$(T - \theta I)\varphi(T) = 0.$$

Since φ is a proper factor of ψ , we see that $\varphi(T) \neq 0$. Let w be a non-zero column of $\varphi(T)$. Then $(T - \theta I)w = 0$, and so w is an eigenvector for T with eigenvalue θ . \square

If ψ_v is the minimal polynomial of T relative to the vector v and $\psi_v(t) = (t - \theta)\varphi(t)$, then $\varphi(T)v$ is an eigenvector for T with eigenvalue θ . When $\dim V$ is small, this provides an effective way of finding eigenvalues.

For example, suppose $\dim V = 2$, and choose a non-zero vector v . If we are very lucky, v is an eigenvector for T . If not, then T^2v is a linear combination of v and Tv , and ψ_v is quadratic. If θ and τ are the zeros of $\psi_v(t)$, then $(T - \theta I)v$ is an eigenvector for T with eigenvalue τ .

2.6 Primary Decomposition

We use the minimal polynomial of an endomorphism to derive a direct sum decomposition of the space on which it acts. We use the following fact: the greatest common divisor of the polynomials $\varphi_1, \dots, \varphi_k$ is 1 if and only if there are polynomials a_1, \dots, a_k such that

$$a_1\varphi_1 + \dots + a_k\varphi_k = 1.$$

2.6.1 Theorem. *Let T be an endomorphism of V with minimal polynomial $\psi(t)$. Suppose that $\psi(t) = \prod_{i=1}^r \psi_i(t)$, where the factors ψ_i are pairwise coprime. Set $\varphi_r = \psi/\psi_r$ and let $a_1(t), \dots, a_r(t)$ be polynomials such that $\sum_i a_i(t)\varphi_i(t) = 1$. Then V is the direct sum of T -invariant subspaces U_i , where U_i is the range of the idempotent $a_i(T)\varphi_i(T)$. The minimal polynomial of $T|_{U_i}$ is $\psi_i(t)$, and $U_i = \ker \psi_i(T)$.*

Proof. Then the greatest common divisor of the polynomials φ_i is 1, and so there are polynomials a_i such that

$$a_1\varphi_1 + \dots + a_r\varphi_r = 1. \quad (2.6.1)$$

Define

$$E_i := a_i(T)\varphi_i(T).$$

Then

$$\sum_{i=1}^r E_i = I.$$

If $i \neq j$ then ψ divides $\varphi_i\varphi_j$, whence

$$E_i E_j = 0.$$

Together the last two equations imply that $E_i^2 = E_i$; thus E_i is an idempotent.

Let U_i denote the range of E_i . If $u \in U_i$ then $E_i u = u$ and so

$$Tu = TE_i u = E_i Tu.$$

Therefore Tu lies in the range of E_i , and therefore U_i is T -invariant. Hence V is a direct sum as described.

Next we show that the minimal polynomial of $T|_{U_1}$ is ψ_1 . Suppose p is a polynomial such that $p(T)U_i = 0$. Then

$$0 = p(T)E_1 = p(T)a_1(T)\varphi_1(T)$$

which implies that $pa_1\varphi_1$ is divisible by ψ and consequently that ψ_1 divides pa_1 . Since ψ_1 divides each of a_2, \dots, a_r , it follows from (2.6.1) that a_1 and ψ_1 are coprime. Hence ψ_1 divides p , and we conclude that ψ_1 is the minimal polynomial of $T \upharpoonright U_1$. Setting 1 equal to i , the general result follows. \square

Remark: If T has minimal polynomial $\psi(t)$, the ring of all polynomials in T is isomorphic to the quotient ring $\mathbb{F}[t]/(\psi(t))$. The preceding theory is a reflection of the structure theory of this ring.

We use the primary decomposition theorem to prove the following fundamental result.

2.6.2 Theorem. *Let T be an endomorphism of the vector space V over the field \mathbb{F} , where \mathbb{F} is algebraically closed. Then there is a diagonalizable endomorphism S and a nilpotent endomorphism N such that S and N are both polynomials in T and $T = S + N$.*

Proof. Let ψ be the minimal polynomial of T . Since \mathbb{F} is algebraically closed, we may write ψ as

$$\psi(t) = \prod_i (t - \theta_i)^{m_i}.$$

Define ψ_i by

$$\psi_i(t) = \frac{\psi(t)}{(t - \theta_i)^{m_i}}.$$

Let E_i denote $\ker \psi_i(T)$. The polynomials ψ_i are coprime (as a set) and so by the primary decomposition theorem, the E_i are pairwise orthogonal idempotents summing to I . Further each E_i is polynomial in T .

Define S by

$$S = \sum_i \theta_i E_i.$$

If $x \in \ker(T - \theta_i)^{m_i}$ then

$$(T - S)^{m_i} x = (T - \theta_i I)^{m_i} x = 0,$$

from which it follows that $T - S$ is nilpotent. As E_i is a polynomial in T , we see that S is too. \square

2.7 The Degree of the Minimal Polynomial

We have seen that the minimal polynomial of an endomorphism T of V is equal to the least common multiple of the minimal polynomials ψ_v , where v runs over the vectors of a basis V . Fortunately something more concrete is true.

2.7.1 Theorem. *If T is an endomorphism of \mathbb{F}^n , then there is a vector x such that the minimal polynomial of T relative to x is the minimal polynomial of T .*

Proof. Assume first that the minimal polynomial ψ of T equals $p(t)^m$, where p is irreducible. Then $p(T)^m = 0$ but $p(T)^{m-1} \neq 0$. Choose a vector x such that $p(T)^{m-1}x \neq 0$. If ϕ is monic and $\phi(T)x = 0$ then ϕ must divide ψ . If ϕ divides p^{m-1} then $\phi(T)x \neq 0$. Consequently $\phi = p^m$.

Now suppose that the minimal polynomial of T has the coprime factorization $\psi_1\psi_2$ and that U_1 and U_2 are the summands of the corresponding direct sum decomposition of \mathbb{F}^n . Let E_1 and E_2 be the associated idempotents. Suppose that x_i is a vector in U_i such that the minimal polynomial of T relative to x_i is ψ_i . If ϕ is monic and

$$\phi(T)(x_1 + x_2) = 0$$

then

$$0 = E_1\phi(T)(x_1 + x_2) = \phi(T)E_1(x_1 + x_2) = \phi(T)x_1.$$

This implies that ψ_2 divides ϕ and a similar argument shows that ψ_1 divides it. So ψ divides ϕ and $x_1 + x_2$ is the vector we need.

An easy induction argument based on the last two paragraphs yields that there is always a vector x such that the minimal polynomial of T is the minimal polynomial of T relative to x .

If the field we are working with is infinite, there is an alternative proof. First, the set of relative minimal polynomials ψ_v is finite, since they are all monic divisors of ψ . Suppose ψ_1, \dots, ψ_r is a list of all the possibilities, and let V_i be the set of vectors v such that $\psi_i(T)v = 0$. Then V_i is a subspace of V and the union of the spaces V_i is V itself. But a vector space over an infinite field cannot be the union of a finite number of proper subspaces, hence $V_i = V$ for some i and ψ_i is the minimal polynomial of T .

2.7.2 Corollary. *If $\dim V = n$ and $T \in \text{End}(V)$, then the degree of the minimal polynomial of T is at most n .*

Proof. If $T \in \text{End}(V)$ has minimal polynomial $\psi(t)$, then there is a vector v in V such that $\psi(t)$ is the minimal polynomial of T relative to v . Hence, if ψ has degree d , the vectors v, Av_1, \dots, Av_{d-1} are linearly independent. Therefore $\dim(V) \geq d$. \square

- (1) Let T be an endomorphism of \mathbb{F}^n and let x_1, \dots, x_n be a basis for \mathbb{F}^n . If ψ_i denotes the minimal polynomial of T relative to x_i , show that the minimal polynomial of T is the least common multiple of ψ_1, \dots, ψ_n .
- (2) Prove that a vector space over an infinite field cannot be the union of a finite number of proper subspaces.

2.8 Root Spaces

We consider primary decomposition when the field of scalars is algebraically closed. In this case, if T is a linear operator on V with minimal polynomial $\psi(t)$, then $\psi(t)$ has the coprime factorization

$$\psi(t) = \prod_{i=1}^k (t - \theta_i)^{m_i},$$

where $\theta_1, \dots, \theta_k$ are the distinct zeros of ψ . It follows that V is the direct sum of the subspaces

$$\ker(T - \theta_i)^{m_i}.$$

We call these subspaces the *root spaces* of T .

If $v \in V$ and $(T - \theta I)^m v = 0$, then the minimal polynomial of T relative to v divides $(t - \theta)^m$. We say that v is a *root vector* for T if its minimal polynomial relative to T has the form $(t - \theta)^r$, for some integer r . If $(T - \theta I)^r v = 0$ and $v \neq 0$, then θ is an eigenvalue of T .

Since V is the direct sum of the root spaces of T , we have the following fundamental result.

2.8.1 Theorem. *Let V be a finite-dimensional vector space over an algebraically closed field. If T is a linear operator on V , then V has a basis consisting of root vectors of T .* \square

The dimension of the root space of an eigenvalue θ of T is called its *algebraic multiplicity*. (The dimension of $\ker(T - \theta I)$ is the *geometric multiplicity* of the eigenvalue.)

2.8.2 Lemma. *Let T be a linear operator on V and let v_1, \dots, v_n be non-zero root vectors. If the respective eigenvalues of these vectors are distinct, then they are linearly independent.*

Proof. Assume $\dim(V) = n$. Suppose that we have scalars a_1, \dots, a_k , not all zero, such that

$$\sum_{i=1}^k a_i v_i = 0. \quad (2.8.1)$$

We prove by induction on k that $a_1 = \dots = a_k = 0$. When $k = 1$, this claim is trivial. Assume $k > 1$. If we apply $(T - \theta_k I)^n$ to both sides of the above expression we get

$$a_1 (T - \theta_k I)^n v_1 + \dots + a_{k-1} (T - \theta_k I)^n v_{k-1} = 0. \quad (2.8.2)$$

Since none of v_1, \dots, v_{k-1} lie in the root space belonging to θ_k , none of the $k-1$ terms in this sum is zero. Since each root space is T -invariant, $(T - \theta_i I)^n v_i$ is therefore a non-zero root vector in the root space containing v_i . So by induction, (2.8.2) implies that $a_1 = \dots = a_{k-1} = 0$. From (2.8.1) it follows that $a_k = 0$ too, and we conclude that v_1, \dots, v_k are linearly independent. \square

(1) Let T be a linear operator on V with an eigenvalue θ . Show that all root vectors belonging to θ are eigenvectors if and only if

$$\ker(T - \theta I) \cap \text{range}(T - \theta I) = \{0\}.$$

2.9 Examples of Root Spaces

We give three examples of root spaces.

Suppose $\dim V = n$ and e_1, \dots, e_n is a basis for V . Then there is a linear operator T on V such that

$$T(e_i) = \begin{cases} e_{i+1}, & \text{if } i < n; \\ 0, & \text{if } i = n. \end{cases}$$

Thus, if $r < n$ then $T^r(e_1) = e_{1+r}$ and $T^n e_1 = 0$. In this case V is the root space belonging to the eigenvalue 0.

Let V be $C^\infty(\mathbb{R})$ and let D be differentiation. Then $\ker D^r$ is the space of polynomials of degree less than r . With some work, we can determine $\ker(D - \lambda I)^r$. First we define a linear operator M_λ on V by

$$M_\lambda(f) := e^{\lambda t} f.$$

We claim that $D - \lambda I = M_\lambda D M_{-\lambda}$. (So D and $D - \lambda I$ are similar.)

To verify this we compute

$$\begin{aligned} DM_{-\lambda}(f) &= \frac{d}{dt} e^{-\lambda t} f(t) \\ &= -\lambda e^{-\lambda t} f(t) + e^{-\lambda t} f'(t) \\ &= e^{-\lambda t} (-\lambda f(t) + D(f(t))) \\ &= M_{-\lambda}(D - \lambda I)f. \end{aligned}$$

Since $M_\lambda^{-1} = M_{-\lambda}$, it follows that for all f in V ,

$$(M_\lambda D M_{-\lambda})(f) = (D - \lambda I)(f),$$

which is what we claimed.

Now we determine $\ker(D - \lambda I)^r$. We have

$$(D - \lambda I)^r = M_\lambda D^r M_{-\lambda}$$

and therefore $(D - \lambda I)^r(g) = 0$ if and only if

$$M_\lambda D^r M_{-\lambda}(g) = 0.$$

Since M_λ is invertible this holds if and only if

$$D^r M_{-\lambda}(g) = 0.$$

Accordingly $\ker(D - \lambda I)^r$ consists of the functions $g(t)$ such that $e^{-\lambda t} g(t)$ is a polynomial of degree less than r . Therefore $\ker(D - \lambda I)^r$ consists of the functions $e^{\lambda t} p(t)$ where $p(t)$ is a polynomial of degree less than r .

Let $V = \mathbb{C}^{\mathbb{N}}$ and let S be the left shift on V . Define a linear operator M_λ by

$$M_\lambda(a_0, a_1, a_2, \dots) := (a_0, \lambda a_1, \lambda^2 a_2, \dots).$$

If $\lambda \neq 0$, then $M_\lambda^{-1} = M_{\lambda^{-1}}$ and

$$S - \lambda I = M_\lambda(S - I)M_\lambda^{-1}.$$

We can show that $\ker(S - I)^r$ consists of the sequences

$$(p(0), p(1), p(2), \dots)$$

where p is a polynomial of degree less than r , and hence we can show that $\ker(S - \lambda I)^r$ consists of the sequences

$$(p(0), \lambda p(1), \lambda^2 p(2), \dots)$$

where p is again a polynomial of degree less than r .

The kernel of S^r consists of the sequences $(a_i)_{i \geq 0}$ such that $a_i = 0$ if $i > r$.

2.10 Differential Equations

We begin with two technical results. In this section V is a vector space over \mathbb{C} .

2.10.1 Lemma. *Let $T : V \rightarrow V$ be linear and suppose that if $\lambda \in \mathbb{C}$, then $\dim(\ker(T - \lambda I)) \leq 1$. If $p(t)$ is a polynomial of degree n , then $\dim(\ker p(T)) \leq n$.*

Proof. We prove the result by induction on the degree of $p(t)$. If $n = 1$, there is nothing to prove. Assume $n > 1$.

Suppose θ is a zero of $p(t)$. Then

$$p(t) = (t - \theta)q(t),$$

where q is a polynomial of degree $n - 1$. By induction on n , we see that $U = \ker q(T)$ has dimension at most $n - 1$.

Now $\ker p(T)$ consists of all vectors v such that $q(T)v$ lies in $\ker(T - \theta I)$. Hence $q(T)$ maps $\ker p(T)$ into $\ker(T - \theta I)$. Let S denote the restriction of $q(T)$ to $\ker p(T)$. Then by the dimension theorem,

$$\dim(\ker p(T)) = \dim \ker(S) + \text{rk}(S) \leq \dim(\ker(q(T))) + 1 \leq n. \quad \square$$

The hypotheses of this lemma hold when $V = C^\infty(\mathbb{R})$ and T is differentiation, or when $V = \mathbb{C}^{\mathbb{N}}$ and T is the left shift.

2.10.2 Theorem. *Let T be a linear operator on V and let $p(t)$ be a polynomial whose zeroes are $\theta_1, \dots, \theta_k$, with respective multiplicities v_1, \dots, v_k . If $\ker p(T)$ has finite dimension, it has a basis consisting of root vectors of T ; the eigenvalues of these root vectors are the zeros of $p(t)$ and the index of the root vectors with eigenvalue θ_i is at most v_i .*

Proof. Suppose $K := \ker p(T)$. If $u \in K$, then

$$p(T)Tu = Tp(T)u = 0$$

and therefore K is T -invariant. Hence K is spanned by root vectors of the restriction of T to K , and these are root vectors of T . Suppose z is a root vector of T with eigenvalue θ and index m . Then

$$(T - \theta I)^m z = 0, \quad p(T)z = 0.$$

Therefore the minimal polynomial of T relative to z divides $(t - \theta)^m$ and $p(t)$, and thus it divides $(t - \theta)^v$, where v is the multiplicity of θ as a zero of $p(t)$. \square

Let $V = C^\infty(\mathbb{R})$ and let D be differentiation. if

$$p(t) := t^n + a_1 t^{n-1} + \cdots + a_n,$$

then the set of solutions to the differential equation

$$D^n f + a_1 D^{n-1} f + \cdots + a_n f = 0$$

is the kernel of $p(D)$. By Lemma 2.10.1 we see that $\ker p(D)$ has finite dimension and so by Theorem 2.10.2, it follows that $\ker p(D)$ is spanned by root vectors of D whose eigenvalues are zeros of $p(t)$.

We want to find all solutions to

$$D^2 f + 3Df + 2f = 0.$$

The solution set of this equation is $\ker p(D)$, where

$$p(t) := t^2 + 3t + 2 = (t + 1)(t + 2)$$

From our work above, this subspace has a basis consisting of root vectors for D . Since the zeros of $p(t)$ are simple we only need root vectors of index one, that is, we only need eigenvectors. Hence the functions

$$e^{-t}, e^{-2t}$$

form a basis for the solution space of this differential equation and therefore every solution can be written as

$$Ae^{-t} + Be^{-2t},$$

for some scalars A and B .

Suppose we want all solutions of

$$D^2 f + 2Df + f = 0$$

Here

$$p(t) = (t + 1)^2,$$

whence we see that $\ker p(D)$ is spanned by root vectors with eigenvalue -1 and index at most two. Therefore it is spanned by

$$e^{-t}, te^{-t};$$

the solutions all have the form

$$(A + Bt)e^{-t}$$

for some scalars A and B .

2.11 Linear Recurrence Equations

The Fibonacci sequence $\varphi = (f_n)_{n \geq 0}$ is defined by the recurrence

$$f_{n+1} = f_n + f_{n-1} \tag{2.11.1}$$

and the initial conditions $f_0 = f_1 = 1$. We want to find an explicit expression for the terms of this sequence.

Let S denote the left shift on $\mathbb{C}^{\mathbb{N}}$. Then we may rewrite (2.11.1) as

$$S^2 \varphi = S\varphi + \varphi;$$

this suggests we should study $\ker p(S)$, where

$$p(t) = t^2 - t - 1.$$

The zeros of $p(t)$ are

$$\frac{1 \pm \sqrt{5}}{2};$$

denote these by θ and τ , where $\theta > \tau$. It follows from Theorem 2.10.2 that $\ker p(S)$ is spanned by root vectors for θ and τ with index at most one, hence by eigenvectors.

The eigenvector for S with eigenvalue a is the geometric series

$$(1, a, a^2, \dots)$$

and therefore there are constants a and b such that

$$f_n = a\theta^n + b\tau^n.$$

Setting $n = 0$ and $n = 1$ here gives two equations in the unknowns a and b :

$$1 = a + b, \quad 1 = a\theta + b\tau.$$

We can rewrite the second equation as

$$1 = \frac{a+b}{2} + \frac{a-b}{2}\sqrt{5};$$

since $a + b = 1$ this implies that

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

Therefore

$$a = \frac{1 + \sqrt{5}}{2\sqrt{5}} = \frac{\theta}{\sqrt{5}}$$

and

$$b = \frac{-1 + \sqrt{5}}{2\sqrt{5}} = -\frac{\tau}{\sqrt{5}}$$

We conclude that

$$f_n = \frac{1}{\sqrt{5}}(\theta^{n+1} - \tau^{n+1}).$$

2.12 Diagonalizability

A matrix A is *diagonalizable* if there is a diagonal matrix D and an invertible matrix L such that $A = LDL^{-1}$, that is, A is similar to a diagonal matrix. If $A = LDL^{-1}$ then $A^k = LD^kL^{-1}$, and so computing A^k can be reduced to the simpler task of computing D^k . More generally, it is often possible to reduce questions about diagonalizable matrices to questions about diagonal matrices (which are often trivial).

2.12.1 Theorem. For an $n \times n$ matrix A over an algebraically closed field \mathbb{F} , the following are equivalent:

- (a) A is diagonalizable.
- (b) \mathbb{F}^n has a basis that consists of eigenvectors of A .
- (c) The minimal polynomial of A has no repeated factors.

Proof. If two matrices are similar, their minimal polynomials are equal, and so (a) implies (c).

If the minimal polynomial has no repeated factors then there are no root vectors of index greater than one, and thus it follows that \mathbb{F}^n has a basis formed from eigenvectors of A .

Finally, suppose that the columns of L are a basis consisting of eigenvectors. Then each column of AL is a scalar multiple of the corresponding column of L , and therefore there is a diagonal matrix D such that $AL = LD$. Since L must be invertible, (a) follows. \square

If \mathbb{F} is not algebraically closed (or close to it, like \mathbb{R}), then diagonalizability is not usually a useful concept.

Chapter 3

Frobenius Normal Form

We derive some properties of matrices from the theory we have established, and then develop the theory of the Frobenius normal form.

3.1 Companion Matrices

Let T be an endomorphism of the finite-dimensional vector space V . One of the best ways to study T is to find T -invariant subspaces of V , and cyclic subspaces are the most accessible of these.

The dimension of the subspace U generated by a vector v is the least integer k such that $T^k v$ lies in the span of the vectors

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

and this set of vectors forms a natural basis for U . Let v_i denote $T^i v$. Then there are scalars a_1, \dots, a_k such that

$$Tv_{d-1} = -a_d v_0 - \dots - a_1 v_{d-1}. \quad (3.1.1)$$

If $i < d - 1$, then $Tv_i = v_{i+1}$ and therefore the matrix representing the action of T on U , relative to the ordered basis v_0, \dots, v_{d-1} , has the form

$$\begin{pmatrix} 0 & 0 & \cdots & 0 & -a_d \\ 1 & 0 & & 0 & -a_{d-1} \\ 0 & 1 & & 0 & -a_{d-2} \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & -a_1 \end{pmatrix} \quad (3.1.2)$$

We call this matrix the *companion matrix* of the polynomial

$$p(t) = t^d + a_1 t^{d-1} + \cdots + a_d.$$

(We will also refer to this as the *right* companion matrix of p ; we will meet other flavours as we proceed.) Since $v_i = T^i v_0$, from (3.1.1) we find that

$$p(T)v_0 = (T^d + a_1 T^{d-1} + \cdots + a_d I)v_0 = 0.$$

Thus $p(t)$ is the minimal polynomial of T relative to v .

We now consider a matrix view of the previous material. Suppose $v \in \mathbb{F}^n$ and $A \in \text{Mat}_{n \times n}(\mathbb{F})$. Assume that the A -cyclic subspace generated by u has dimension d and let the matrix R be given by

$$R := (u \quad Au \quad \cdots \quad A^{d-1}u).$$

Thus $\text{col}(R)$ is the A -cyclic subspace generated by u . If $\psi(t)$ is the minimal polynomial of A relative to u and C is the companion matrix of ψ , then

$$AR = (Au \quad A^2u \quad \cdots \quad A^d u) = RC.$$

There is a third view, which is also quite important. Suppose ψ is a polynomial of degree d over \mathbb{F} and let V_ψ be the vector space of polynomials over \mathbb{F} modulo ψ . This vector space is usually denoted by $\mathbb{F}[z]/(\psi(z))$; its elements are equivalence class of polynomials, where polynomials f and g are equivalent if and only if $f - g$ is divisible by ψ . Each equivalence class contains a unique polynomial of degree less than d , and these are the natural representatives of the equivalence classes.

The powers

$$1, z, \dots, z^{d-1}$$

provide one basis for V_ψ . Multiplication by z is an endomorphism of V_ψ , and the matrix representing multiplication by z relative to this basis is easily seen to be the companion matrix of ψ .

1. Let $p(z)$ be a polynomial of degree k as above and let C_p denote its companion matrix. If f is a polynomial of degree less than k , let \hat{f} be the coordinate vector of f relative to the standard basis $1, x, \dots, x^{k-1}$. Use the fact that $f(z)z^i$ and $z^i f(z)$ have the same remainder modulo p to prove that

$$f(C_p) = (\hat{f} \quad C_p \hat{f} \quad \cdots \quad (C_p)^{k-1} \hat{f}).$$

Deduce that $f(C_p)\hat{g} = g(C_p)\hat{f}$.

2. If C_p is a companion matrix of order $n \times n$, show that $\text{rk}(C_p - \theta I) \geq n - 1$, for any element θ of \mathbb{F} . Deduce that the geometric multiplicity of any eigenvalue is at most 1. (This implies that C is diagonalizable if and only if the zeros of p are all simple.)
3. Let U be the subspace spanned by the vectors $T^r u$, where $r \geq 0$. If $Su \in U$, show that there is a polynomial p such that $Su = p(T)u$. Hence deduce that if U is S -invariant and $ST = TS$, then $S|_U$ is a polynomial in $T|_U$.

3.2 Transposes

We introduce a second basis for V_ψ . If

$$\psi(z) = t^d + a_1 t^{d-1} + \cdots + a_d,$$

define polynomials ψ_1, \dots, ψ_d by

$$\psi_i(z) := t^{d-i} + a_1 t^{d-i-1} + \cdots + a_{d-i}.$$

These polynomials can also be defined by the initial condition $\psi_d(z) = 1$ and the backwards recurrence

$$\psi_{i-1}(z) = z\psi_i(z) + a_{d-i+1}. \quad (3.2.1)$$

As a third alternative, we can view $\psi_i(z)$ as the polynomial part of the rational function $z^{-i}\psi(z)$. Since $\psi_i(z)$ is monic of degree $d - i$, we see that these polynomials form a basis for V_ψ , sometimes called the *control basis*.

Suppose $v \in V$ and T is an endomorphism of V with minimal polynomial $\psi(z)$ relative to v . Then the vectors

$$\psi_1(T)v, \dots, \psi_d(T)v$$

form a basis for the T -cyclic subspace U generated by v . It follows from (??) that

$$T\psi_i(T)v = \begin{cases} -a_d v, & \text{if } i = 1; \\ \psi_{i-1}(T)v - a_{d+1-i} v, & \text{if } 2 \leq i \leq d. \end{cases}$$

From this we see that the matrix representing T with respect to the control basis is C_ψ^T , the transpose of the companion matrix of ψ .

It follows that since C and C^T both represent the action of T on U , they are similar. We can say more about this.

Let Q be the $d \times d$ matrix given by

$$Q = \begin{pmatrix} a_{d-1} & a_{d-2} & \cdots & a_1 & 1 \\ a_{d-2} & a_{d-3} & & 1 & 0 \\ \vdots & & & & \vdots \\ a_1 & 1 & & & 0 \\ 1 & 0 & \cdots & & 0 \end{pmatrix}.$$

Note that Q is symmetric and invertible. Choose an ordered basis α for V and let A be the matrix representing T in its action of V . Let R be the matrix with the coordinate vectors $[T^i v]_\alpha$ as its columns and let S be the matrix with the vectors $[p_i(T)v]_\alpha$ as its columns. Then $S = RQ$ and

$$AR = RC, \quad AS = SC^T.$$

Hence $ARQ = AS = SC^T = RQC^T$ and therefore $RCQ = RQC^T$. Since the columns of R are linearly independent, it follows that $QC^T = CQ$.

It will follow from what we have just proved that if A is an $n \times n$ matrix, there is a symmetric invertible matrix Q such that $A^T = Q^{-1}AQ$.

3.3 Eigenvectors for Companion Matrices

We give explicit formulas for the left and right eigenvectors of a companion matrix. We use e_1, \dots, e_d to denote the standard basis vectors of \mathbb{F}^d , as customary.

3.3.1 Lemma. *Let $\psi(z)$ be a polynomial of degree d and let C be its companion matrix. Then*

$$(1 \ z \ \cdots \ z^{d-1})C = z(1 \ z \ \cdots \ z^{d-1}) - \psi(z)e_d^T.$$

Proof. Suppose

$$\psi(z) = t^d + a_1 t^{d-1} + \cdots + a_d.$$

If $i < d$, the i -th entry of

$$(1 \ z \ \cdots \ z^{d-1})C$$

is z^{i+1} ; while the d -th entry is

$$-(a_1 + a_2 z + \cdots + a_d z^{d-1}) = z^d - \psi(z).$$

The lemma follows at once from this. □

If, in the above lemma, we take z to be a zero θ of ψ , then it follows that

$$(1 \quad \theta \quad \dots \quad \theta^{d-1})$$

is a left eigenvector of C with eigenvalue θ .

Our next lemma will provide right eigenvectors. Let ψ_1, \dots, ψ_d denote the control basis for V_ψ .

3.3.2 Lemma. *Let $\psi(z)$ be a polynomial of degree d and let C be its companion matrix. Then*

$$C \begin{pmatrix} \psi_1 \\ \vdots \\ \psi_d \end{pmatrix} = z \begin{pmatrix} \psi_1 \\ \vdots \\ \psi_d \end{pmatrix} - \psi(z)e_1.$$

Proof. This is again routine; we leave it as an exercise. \square

These two lemmas provides right and left eigenvectors for C , one for each zero θ of $\psi(z)$. If $\psi(z)$ has d distinct zeros, we obtain d distinct left eigenvectors for C . Since the eigenvalues are distinct, these eigenvectors are linearly independent.

If we are working over \mathbb{R} or \mathbb{C} , we can say something useful when $\psi(z)$ has zeros with multiplicity greater than 1. The idea is to differentiate both sides of the identity in Lemma 3.3.2. Define

$$\Psi(z) := \begin{pmatrix} \psi_1 \\ \vdots \\ \psi_d \end{pmatrix}$$

and let $\Psi^{(r)}(z)$ denote the r -th derivative of $\Psi(z)$. Then

$$C\Psi^{(r)}(z) = z\Psi^{(r)}(z) + r\Psi^{(r-1)}(z) - \psi^{(r)}(z)e_1.$$

If θ is a zero of ψ with multiplicity m and $r < m$, then $\psi^{(r)}(\theta) = 0$ and therefore

$$(C - \theta I)^r \Psi^{(r)}(\theta) = r! \Psi(\theta).$$

Note the since the polynomials ψ_i form a basis for the polynomials of degree less than d , they cannot all be zero at θ ; therefore $\Psi(\theta) \neq 0$. It follows that the vectors $\Psi^{(r)}(\theta)$ are a basis for the root space associated with θ . (Exercise: show that these vectors are linearly independent.)

1. By expanding the expression

$$(1 \quad w \quad \cdots \quad w^{d-1}) C \begin{pmatrix} \psi_1(z) \\ \vdots \\ \psi_d(z) \end{pmatrix}$$

in two different ways, derive the identity

$$(w - z) \sum_i w^i \psi_i(z) = \psi(w) - \psi(z).$$

(If we take w and z to be zeros of ψ , this gives the orthogonality relation between the right and left eigenvectors of C .)

2. Let Q be the symmetric matrix from ???. Show that

$$\Psi(z) = Q \begin{pmatrix} 1 \\ z \\ \vdots \\ z^{d-1} \end{pmatrix},$$

and hence deduce that $C^T = Q^{-1} C Q$.

3.4 Inverses of Companion Matrices

Suppose $A \in \text{Mat}_{n \times n}(\mathbb{F})$ and that $u \in \mathbb{F}^n$ that generates an A -cyclic subspace of dimension d . Let ψ be the minimal polynomial of A relative to u . If

$$\psi(t) := t^k + a_1 t^{k-1} + \cdots + a_k,$$

then C is invertible if and only if $a_k \neq 0$. (There are a number of ways to see this. Perhaps the easiest is to note that if we move the last column of C to the first position, the resulting matrix C' is lower triangular with $(C')_{1,1} = -a_k$ and all other diagonal entries equal to 1.) If C is invertible, there is a simple expression for C^{-1} . To describe this, we need a new operation on polynomials.

If q is a polynomial with degree k , let \tilde{q} denote the polynomial $t^k q(t^{-1})$. (This is q with its coefficients reversed.) Note that if A is invertible, then $p(A) = 0$ if and only if

$$A^k \tilde{p}(A^{-1}) = 0.$$

It follows that if p is the minimal polynomial of A , then $a_k^{-1}\tilde{p}$ is the minimal polynomial of A^{-1} .

Let R be the matrix given by

$$R = (u \quad Au \quad \cdots \quad A^{d-1}u).$$

Then $AR = RC$ and $\text{col}(R)$ is the A -cyclic subspace generated by u . If A is invertible, then A^{-1} is a polynomial in A and therefore $\text{col}(R)$ is A^{-1} -invariant. Hence there is a matrix D such that

$$A^{-1}R = RD$$

and $D = C^{-1}$. Now

$$A^{-1}(u \quad Au \quad \cdots \quad A^{d-1}u) = (A^{-1}u \quad u \quad \cdots \quad A^{d-2}u),$$

whence D is a $d \times d$ matrix of the form

$$\begin{pmatrix} \gamma & I_{d-1} \\ c_d & 0 \end{pmatrix}.$$

If we write C in the form

$$C = \begin{pmatrix} 0 & a_d \\ I_{d-1} & \alpha \end{pmatrix},$$

then the equation $DC = I$ implies that

$$I_d = \begin{pmatrix} I_{d-1} & a_d\gamma + \alpha \\ 0 & a_d c_d \end{pmatrix}$$

Consequently we must have

$$c_d = a_d^{-1}, \quad \gamma = -a_d^{-1}\alpha$$

and therefore

$$D = \begin{pmatrix} -a_d^{-1}\alpha & I_{d-1} \\ a_d^{-1} & 0 \end{pmatrix}.$$

This expression for D makes sense if and only if $a_d \neq 0$, because C can be invertible even when A is not. Hence we have proved the following:

3.4.1 Theorem. *Let p be a polynomial with degree k and let C be the companion matrix of p . Then C is invertible if and only if $p(0) \neq 0$. If $p(0) \neq 0$, then*

$$C^{-1} = TDT,$$

where D is the companion matrix of $a_k^{-1}\tilde{p}$ and T is the matrix whose columns are the standard basis vectors in reverse order. \square

By way of example, we have

$$\begin{pmatrix} -c/d & 1 & 0 & 0 \\ -b/d & 0 & 1 & 0 \\ -a/d & 0 & 0 & 1 \\ -1/d & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & 0 & 0 & -d \\ 1 & 0 & 0 & -c \\ 0 & 1 & 0 & -b \\ 0 & 0 & 1 & -a \end{pmatrix} = I.$$

If C is a companion matrix and T is the permutation matrix in the previous theorem, we say that TCT is a *left companion matrix*. Analogously we will call C^T a *bottom companion matrix*. And to round off the list, $TC^T T$ is a *top companion matrix*. All four flavours occur in practice.

3.5 Cycles

Let P be the $n \times n$ matrix such that $Pe_1 = e_n$ and, if $2 \leq i \leq n$ then $Pe_i = e_{i-1}$ and $Pe_n = e_1$. Thus if $n = 5$,

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

We see that $P^n = I$ and P is a companion matrix for the polynomial $t^n - 1$. Further $P^{-1} = P^T$, and therefore P is orthogonal.

Let v_θ be the vector of length n with i -th entry θ^{i-1} . Thus

$$v_\theta = \sum_i \theta^{i-1} e_i$$

and consequently, if $\theta^n = 1$, then

$$Pv_\theta = \sum_i \theta^{i-1} Pe_i = \sum_i \theta^{i-1} e_{i-1} = \theta v_\theta.$$

Therefore the vectors v_θ , as θ runs over the distinct n -th roots of unity, are eigenvectors for P . It is not hard to show that, if n is odd, any real eigenvector of P is a scalar multiple of v_1 .

Now let $A = P + P^T$. Then A is symmetric and

$$Av_\theta = (\theta + \theta^{-1})v_\theta.$$

Therefore the vectors v_θ , as θ runs over the distinct n -th roots of unity, are eigenvectors for A . Note that here the eigenvalues

$$\theta + \theta^{-1} = \theta + \bar{\theta}$$

are real, even though the eigenvectors themselves are complex (unless θ is real).

The eigenvalues of P are roots of unity. Suppose Q is orthogonal and v is an eigenvector for it with eigenvalue θ . Then $Qv = \theta v$, but

$$\|v\| = \|Qv\| = \|\theta v\| = |\theta| \|v\|.$$

It follows that all eigenvalues of an orthogonal matrix lie on the unit circle in the complex plane.

3.6 Circulants and Cyclic Codes

Let P_n be the companion matrix for the polynomial $t^n - 1$. Thus if $n = 5$ then

$$P_5 = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

We see that $P_n e_i = e_{i+1}$ if $i < n$ and $P_n e_n = e_1$. A *circulant matrix* is a matrix which is a polynomial in P_n . This is equivalent to stating that a matrix is a circulant if it is square and each row is a cyclic right shift of the row above it. If the first column of the circulant A is

$$\begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix}$$

then

$$A = \sum_{i=1}^n a_i P^{i-1}.$$

It follows that there is an isomorphism between the vector space of $n \times n$ circulant matrices and the space of polynomials with degree less than n . But this is misleading. Suppose a and b are polynomials with degree less than n , and associated circulants A and B respectively. Then the product AB is a circulant, but

the polynomial belonging to it cannot be ab unless the degree of this product is less than n . In fact the polynomial is the remainder of the product $a(t)b(t)$ on division by $t^n - 1$. Thus the space of $n \times n$ circulants is isomorphic to the quotient ring $\mathbb{F}[x]/(x^n - 1)$. This isomorphism is an algebra isomorphism. If $\deg(a) < n$, we use C_g to denote the circulant associated with g .

The row space of an $n \times n$ circulant over \mathbb{F} is a *cyclic code* of length n . Suppose f is a polynomial and that g is its greatest common divisor with $t^n - 1$. Then there are polynomials a and b such that

$$a(t)f(t) + b(t)(t^n - 1) = g(t).$$

Hence

$$C_g = C_{af} = C_a C_f$$

and therefore $\text{row}(C_g) \subseteq \text{row}(C_f)$. On the other hand, $f = f_1 g$ and so

$$C_f = C_{f_1} C_g,$$

which implies that $\text{row}(C_f) \subseteq \text{row}(C_g)$ and hence that $\text{row}(C_f) = \text{row}(C_g)$. This proves that a cyclic code of length n over \mathbb{F} is equal to C_g , for some divisor g of $t^n - 1$.

One of the most important parameters of a code is its dimension. Thus we would like to determine $\text{rk}(C_g)$. If g has degree d , then the submatrix formed by the intersection of the first $n - d$ columns and last $n - d$ rows of C_g is the identity matrix I_{n-d} . Therefore

$$\text{rk}(C_g) \geq n - d.$$

Suppose $a(t)$ is a polynomial of degree less than n , and let $[a]$ denote its coordinate vector with respect to the ordered basis $1, t, \dots, t^{n-1}$. If $C_g[a] = 0$, then

$$0 = C_g[a] = P^r C_g[a] = C_g P^r [a]$$

for all r and consequently

$$C_g C_a = 0.$$

Equivalently, $C_g C_a = 0$ if and only if $C_g C_a e_1 = 0$. Now $C_g C_a = 0$ if and only if $t^n - 1$ divides $g(t)a(t)$, and accordingly the null space of C_g consists of the vectors $[a]$ such that $(t^n - 1)/g(t)$ divides $a(t)$. If we set $h(t)$ equal to $(t^n - 1)/g(t)$, then the null space of C_g is the column space of C_h . The dimensions of the row and column spaces of C_h are equal, and therefore

$$\text{rk}(C_h) \geq n - (n - d) = d.$$

So by the rank+nullity theorem,

$$\text{rk}(C_g) + \text{rk}(C_h) = n,$$

which forces us to conclude that $\text{rk}(C_g) = n - d$.

If $C_{g^2}[a] = 0$ then $C_{g^2}C_a = 0$ and so $t^n - 1$ divides $g(t)^2a(t)$. If $t^n - 1$ has no repeated factors, then $t^n - 1$ divides $g(t)^2a(t)$ if and only if it divides $g(t)a(t)$. In this case it follows that C_g is diagonalizable. If $x^n - 1 = p(t)^2q(t)$, then $C_{p^2q}^2 = 0$ and C_{pq} is nilpotent, and not diagonalizable.

However $x^n - 1$ has a repeated factor if and only if its gcd with its derivative nx^{n-1} is not constant, in other words, if and only if n is not divisible by the characteristic of \mathbb{F} . In particular, if the characteristic of \mathbb{F} does not divide n , then \mathbb{F}^n is the direct sum of $\ker(C_g)$ and $\text{col}(C_g)$.

Let \mathbb{E} be an extension field of \mathbb{F} in which $t^n - 1$ splits into linear factors. If n and the characteristic of \mathbb{F} are coprime, these factors are all distinct. It follows that each divisor g of $t^n - 1$ is determined by the set of n -th roots of 1 on which it vanishes. Let v_θ be the vector of length n with i -th entry equal to θ^{i-1} . Then if $\theta^n = 1$,

$$C_g v_\theta = g(\theta^{-1})$$

and so $\text{row}(C_g)$ consists of the vectors x^T such that

$$x^T v_\theta = 0$$

whenever θ^{-1} is a zero of g .

3.7 Frobenius Normal Form

A square matrix C is in *Frobenius normal form* if

- (a) It is block-diagonal, with diagonal blocks C_1, \dots, C_m .
- (b) Each diagonal block is the companion matrix of a polynomial $\psi_i(t)$.
- (c) For $i = 1, \dots, m - 1$, the polynomial ψ_{i+1} divides ψ_i .

Thus the Frobenius normal form can be specified by giving the sequence of polynomials ψ_i .

We want to prove that two matrices over a field are similar if and only if they have the same Frobenius normal form. We require two preliminary results.

3.7.1 Lemma. *Then if N_1 and N_2 are similar, they are equal.*

Proof. Suppose that

$$N_1 := \begin{pmatrix} L_1 & 0 \\ 0 & D_1 \end{pmatrix}, \quad N_2 := \begin{pmatrix} L_2 & 0 \\ 0 & D_2 \end{pmatrix}$$

are both in Frobenius normal form. Then $p(N) = 0$ if and only if $p(L_1) = 0$ and $p(D_1) = 0$. Hence the minimal polynomial of N_1 is the minimal polynomial of L_1 . Since N_1 and N_2 are similar, they have the same minimal polynomial, and this is also the minimal polynomial of L_2 . Thus L_1 and L_2 have the same minimal polynomial, and as they are companion matrices this implies that they are equal.

Now let ψ_1 be the minimal polynomial of D_1 . Then $\psi_1(N_1)$ and $\psi_1(N_2)$ are similar and thus

$$\begin{pmatrix} \psi_1(L) & 0 \\ 0 & 0 \end{pmatrix} \sim \begin{pmatrix} \psi_1(L) & 0 \\ 0 & \psi_1(D_2) \end{pmatrix}.$$

This implies that $\psi_1(D_2) = 0$ (prove it!) and we conclude that D_1 and D_2 have the same minimal polynomial. An easy induction argument now yields that $D_1 = D_2$. \square

3.7.2 Lemma. *Let A be an $n \times n$ matrix over \mathbb{F} . If there is a non-zero cyclic subspace U of dimension k , then there is a cyclic subspace of \mathbb{F}^n of dimension at least k with an A -invariant complement. If $\dim(U)$ equals the degree of the minimal polynomial of A , then U has an A -invariant complement.*

Proof. Let u be a non-zero vector and suppose that the A -invariant subspace it generates has dimension k . Let U be the $n \times k$ matrix with the vectors

$$u, Au, \dots, A^{k-1}u$$

as its columns. Then $\text{rk}(U) = k$, and by Lemma ??, there is a $n \times k$ matrix V such that $V^T U = I$. Let w denote the last column of V . (Now that we have w , we will ignore V .)

We have

$$w^T A^r A^s u = w^T A^{r+s} u = \begin{cases} 1, & \text{if } r = k-1-s; \\ 0, & \text{if } r < k-1-s. \end{cases}$$

If W is the matrix with columns

$$(A^{k-1})^T w, (A^{k-2})^T w, \dots, A^T w, w$$

then $W^T U$ is a lower triangular matrix with diagonal entries equal to 1. Therefore it is invertible, and therefore $\text{rk}(W) = k$.

Let ℓ be the dimension of the A^T -invariant subspace generated by w . Since $\text{rk}(W) = k$, we see that $k \leq \ell$. If $k \neq \ell$, then repeating the above argument with A^T in place of A and w in place of U , we obtain a cyclic subspace for A with dimension at least ℓ . By repeating both these steps a finite number of times, we reduce to the case where $k = \ell$. Therefore we may assume that $\text{col}(W)$ is A^T -invariant, and so there is a matrix L such that $A^T W = W L^T$. If $W^T x = 0$ then

$$0 = L W^T x = W^T A x;$$

accordingly the null-space K of W^T is A -invariant. Since $W^T U$ is invertible, no non-zero element of $\text{col}(U)$ lies in K . Since $\text{rk}(W) = k$, we see that $\dim K = n - k$ and therefore K is an A -invariant complement to $\text{col}(U)$.

To obtain the last statement of the proof, note that A and A^T have the same minimal polynomial. So if k equals the degree of this polynomial, then $\text{rk}(U) = \text{rk}(W)$. \square

It follows readily from this lemma that every square matrix is similar to a block diagonal matrix, where each block is a companion matrix. We can also use it as follows to verify the existence of the Frobenius normal form.

3.7.3 Theorem. *Every square matrix is similar to a matrix in Frobenius normal form.*

Proof. Let A be an $n \times n$ matrix with minimal polynomial $\psi(t)$ of degree k . By ??, there is a vector u such that ψ is the minimal polynomial of A with respect to u , and therefore u generates a cyclic subspace V of dimension k . By the previous lemma, it follows that this subspace has an A -invariant complement K .

Choose a basis for \mathbb{F}^n consisting of the columns of V followed by a basis for K . Relative to this basis, A is represented by a block-diagonal matrix

$$\begin{pmatrix} L & 0 \\ 0 & B \end{pmatrix},$$

where L is the companion matrix of the minimal polynomial of A . The minimal polynomial of B divides the minimal polynomial of A . By induction on n we see that B is similar to a matrix in Frobenius normal form; stacking L on top of this produces a matrix in Frobenius normal form that is similar to A . \square

We can use Lemma 3.7.2 to compute the minimal polynomial of a square matrix. First compute a block-diagonal matrix similar to A , with companion matrices as its blocks. The least common multiple of the polynomials associated to these companion matrices is the minimal polynomial of A .

3.8 Applications

We use $\mathcal{C}(M)$ to denote the *commutant* of M , that is, the set of matrices that commute with M . This is a subspace that contains all polynomials in M .

3.8.1 Theorem. *Let A and B be square matrices. If $\mathcal{C}(A) \subseteq \mathcal{C}(B)$, then B is a polynomial in A .*

Proof. Assume A is $n \times n$. We can decompose \mathbb{F}^n as the direct sum of A -invariant subspaces V_1, \dots, V_k . For each subspace there is a cyclic vector v_i such that the powers $A^r v_i$ span V_i . If ψ_i is the minimal polynomial of $A|_{V_i}$, then its degree equals $\dim V_i$, and ψ_{i+1} divides ψ_i .

Let P_i denote the projection on V_i . From ??, the projections P_i commute with A . Hence they commute with B and, again by ??, it follows that the subspaces V_i are B -invariant. Therefore $Bv_i \in V_i$ and therefore there is a polynomial g_i such that $Bv_i = g_i(A)v_i$. As $AB = BA$, we have

$$BA^r v_i = A^r Bv_i = A^r g_i(A)v_i = g_i(A)A^r v_i,$$

and therefore $Bv = g_i(A)v$ for all v in V_i .

To complete the proof, we show that $g_i(A)v_i = g_1(A)v_i$. This implies that $B = g_1(A)$.

Let $q_i := \psi_1/\psi_i$. Consider the map that sends $f(A)v_i$ to $q_i(A)f(A)v_1$. If $f(A)v_i = 0$, then ψ_i divides f and so $\psi_1 = \psi_i q_i$ divides $q_i(A)f(A)v_1$. It follows that this is a well-defined linear map from V_i to V_1 . We extend it a linear map X_i from V to V by defining $X_i(v) = 0$ if $v \in V_j$ and $j \neq i$; if $v = f(A)v_i$ then $X_i v = q_i(A)f(A)v_1$.

If $i \neq j$ and $v \in V_j$, then $AX_i v = X_i A v = 0$. Further

$$AX_i f(A)v_i = Aq_i(A)f(A)v_1$$

and

$$X_i A f(A)v_i = q_i(A)A f(A)v_1.$$

Therefore X_i commutes with A , and therefore it commutes with B . Now

$$X_i B v_i = X_i g_i(A) v_1 = g_i(A) q_i(A) v_1$$

and

$$B X_i v_i = B q_i(A) v_1 = q_i(A) B v_1 = q_i(A) g_1(A) v_1.$$

Since X_i and B commute, this implies that

$$(g_i(A) - g_1(A)) q_i(A) v_1 = 0,$$

whence $(g_i - g_1)q_1$ is divisible by $p_1 = q_i p_i$, and so p_i divides $g_i - g_1$. Consequently

$$g_i(A) v_i = g_1(A) v_i,$$

for all i . □

The above proof follows Prasalov.

3.9 Nilpotent Matrices

A linear mapping or a matrix is *nilpotent* if some power of it is zero. The canonical example is

$$N_2 := \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix},$$

whose square is zero. If T is nilpotent then its minimal polynomial is t^k for some k , sometimes called the *index of nilpotency* of T (but not very often, if we can help it). A nilpotent matrix of index 1 is the zero matrix. We note that N_2 is the companion matrix of t^2 . More generally the companion matrix of t^k is a nilpotent matrix with index k , which we will denote by N_k . Note that $N_k e_1 = 0$ and $N_k e_{i+1} = e_i$ when $i \geq 1$.

Nilpotent matrices are interesting and useful, but also a source of difficulties. Since $N_k e_1 = 0$, we see that e_1 is an eigenvector of N_k with eigenvalue 0. Since the minimal polynomial of N_k is t^k , we see that 0 is the only eigenvalue of N_k . Further, since $\text{rk}(N_k) = k - 1$, the eigenspace associated with 0 has dimension 1, and therefore equals the span of e_1 . Consequently eigenvalues and eigenvectors provide very little information about nilpotent matrices.

We have the following structure theorem.

3.9.1 Theorem. *If M is a nilpotent matrix, then it is similar to a block diagonal matrix, where each diagonal block is equal to N_k for some k .*

Proof. The required block diagonal matrix is the Frobenius normal form of M . \square

One corollary of this is that the number of similarity classes of $n \times n$ nilpotent matrices over a field equals the number of vectors of non-negative integers

$$(k_1, \dots, k_n)$$

such that $k_1 \geq k_2 \geq \dots \geq k_n$ and $\sum_i k_i = n$.

3.9.2 Lemma. *Let A be an $n \times n$ matrix over an algebraically closed field with minimal polynomial $\psi(t)$. Then A is similar to a block diagonal matrix with diagonal blocks of the form $\theta I + N_\theta$, where θ runs over the zeros of ψ , and N_θ is nilpotent with index equal to the multiplicity of θ as a zero of $\psi(t)$.*

Proof. By the primary decomposition theorem Theorem 2.6.1, we know that A is similar to a diagonal matrix with diagonal blocks A_θ indexed by the zeros of ψ , such that the minimal polynomial of A_θ is $(t - \theta)^{m_\theta}$, where m_θ is the multiplicity of θ as a zero of $\psi(t)$. Hence $A - \theta I$ is nilpotent, with index m_θ . Thus we may write

$$A_\theta = \theta I + N_\theta,$$

where N_θ is nilpotent, of index m_θ . \square

The corank of $(A - \theta I)^{m_\theta}$ is known as the *algebraic multiplicity* of the eigenvalue θ . This distinguishes it from the *geometric multiplicity*, which is the corank of $A - \theta I$.

We present one application. We wish to determine when a matrix A has a square root, that is, when there is a matrix X such that $X^2 = A$. If $A = LBL^{-1}$ and B has a square root Y , then

$$(LYL^{-1})^2 = LY^2L^{-1} = LBL^{-1} = A.$$

This allows us to use the primary decomposition theory; more precisely, we assume that A is block diagonal with blocks of the form

$$A_\theta = \theta I + N_\theta.$$

It follows that A has a square root if and only if each of its blocks does.

Suppose M is nilpotent. Then $I + M$ has a square root

$$(I + M)^{1/2} = \sum_{r \geq 0} \binom{\frac{1}{2}}{r} M^r.$$

Note that this is a finite sum, since $M^r = 0$ for all but finitely many values of r . It follows that, if $\theta \neq 0$, then

$$\theta I + N_\theta = \theta(I + \theta^{-1}N_\theta)$$

has a square root. Hence we are left with the case where $\theta = 0$, and this our questions reduces to deciding which nilpotent matrices have a square root. If N is nilpotent with index k and $X^2 = N$, then $X^{2k} = 0$ and so X is nilpotent with index $2k$. (This implies that the matrices N_k are not squares.)

Assume now that N is in Frobenius normal form. We claim that the $(k+1) \times (k+1)$ matrix

$$N'_k := \begin{pmatrix} N_k & 0 \\ 0 & 0 \end{pmatrix}$$

has a square root. (You do it!) It follows that N has a square root if and only if its corank is at least as large as the number of non-zero blocks.

3.10 A Similarity Condition

We are given the following two $n \times n$ matrices:

$$\begin{pmatrix} A & 0 \\ 0 & D \end{pmatrix}, \quad \begin{pmatrix} A & B \\ 0 & D \end{pmatrix},$$

where A and D are square. We ask for which matrices B are they similar.

We note that

$$\begin{pmatrix} I & -X \\ 0 & I \end{pmatrix} \begin{pmatrix} A & B \\ 0 & D \end{pmatrix} \begin{pmatrix} I & X \\ 0 & I \end{pmatrix} = \begin{pmatrix} A & AX - XD + B \\ 0 & D \end{pmatrix},$$

and deduce that they are similar if there is a matrix X such that

$$AX - XD = B.$$

We show that this condition is necessary.

Suppose that our two matrices are similar. Then there is an invertible matrix S such that

$$S \begin{pmatrix} A & 0 \\ 0 & D \end{pmatrix} = \begin{pmatrix} A & B \\ 0 & D \end{pmatrix} S.$$

We define linear mappings T_1 and T_2 on the space of $n \times n$ matrices by

$$T_1(Y) := \begin{pmatrix} A & 0 \\ 0 & D \end{pmatrix} Y - Y \begin{pmatrix} A & 0 \\ 0 & D \end{pmatrix}, \quad T_2(Y) := \begin{pmatrix} A & B \\ 0 & D \end{pmatrix} Y - Y \begin{pmatrix} A & 0 \\ 0 & D \end{pmatrix}.$$

We have

$$\begin{aligned} S(T_1(S^{-1}Y)) &= S \begin{pmatrix} A & 0 \\ 0 & D \end{pmatrix} S^{-1}Y - Y \begin{pmatrix} A & 0 \\ 0 & D \end{pmatrix} \\ &= \begin{pmatrix} A & B \\ 0 & D \end{pmatrix} Y - Y \begin{pmatrix} A & 0 \\ 0 & D \end{pmatrix} \\ &= T_2(Y), \end{aligned}$$

and therefore $\ker(T_1)$ and $\ker(T_2)$ have the same dimension.

Let Y be the matrix

$$Y = \begin{pmatrix} Y_{1,1} & Y_{1,2} \\ Y_{2,1} & Y_{2,2} \end{pmatrix},$$

where the partitioning is compatible with the partitioning of the other matrices above. Then

$$T_1(Y) = \begin{pmatrix} AY_{1,1} - Y_{1,1}A & AY_{1,2} - Y_{1,2}D \\ DY_{2,1} - Y_{2,1}B & DY_{2,2} - Y_{2,2}D \end{pmatrix}$$

and

$$T_2(Y) = \begin{pmatrix} AY_{1,1} - Y_{1,1}A + BY_{2,1} & AY_{1,2} - Y_{1,2}D + BY_{2,2} \\ DY_{2,1} - Y_{2,1}D & DY_{2,2} - Y_{2,2}D \end{pmatrix}.$$

We note that

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

lies in $\ker(T_2)$ if and only if $AY_{1,2} - Y_{1,2}D - B = 0$, and we can prove our claim by showing that there is a matrix of this form in $\ker(T_2)$.

Let \mathcal{T}_i denote the restriction to $\ker T_i$ of the linear map

$$\begin{pmatrix} Y_{1,1} & Y_{1,2} \\ Y_{2,1} & Y_{2,2} \end{pmatrix} \mapsto (Y_{2,1} \quad Y_{2,2}).$$

We will prove that \mathcal{T}_1 and \mathcal{T}_2 have the same image.

From the expressions for $T_1(Y)$ and $T_2(Y)$, we see that $\ker \mathcal{T}_1 = \ker \mathcal{T}_2$. Further

$$\operatorname{im} \mathcal{T}_1 = \{(Y_{2,1} \ Y_{2,2}) : DY_{2,1} - Y_{2,1}A = 0, DY_{2,2} - Y_{2,2}D = 0\}$$

and $\operatorname{im} \mathcal{T}_2$ consists of the elements of $\ker \mathcal{T}_1$ for which there are matrices $Y_{1,1}$ and $Y_{1,2}$ such that

$$CY_{2,1} = Y_{1,1}A - AY_{1,1}, \quad CY_{2,2} = Y_{1,2}D - AY_{1,2}.$$

It follows that $\operatorname{im} \mathcal{T}_2 \subseteq \operatorname{im} \mathcal{T}_1$. Now

$$\begin{aligned} \dim((\mathcal{T}_1)) + \operatorname{rk}(\mathcal{T}_1) &= \dim(\ker T_1) \\ \dim((\mathcal{T}_2)) + \operatorname{rk}(\mathcal{T}_2) &= \dim(\ker T_2). \end{aligned}$$

Since \mathcal{T}_1 and \mathcal{T}_2 have the same corank and since T_1 and T_2 have the same corank, it follows that \mathcal{T}_1 and \mathcal{T}_2 have the same rank.

Finally, it easy to verify that

$$\begin{pmatrix} 0 & 0 \\ 0 & -I \end{pmatrix} \in \ker T_1$$

whence

$$(0 \ -I) \in \operatorname{im} \mathcal{T}_1 = \operatorname{im} \mathcal{T}_2$$

and accordingly there is a matrix in $\ker T_2$ of the form

$$\begin{pmatrix} Y_{1,1} & Y_{1,2} \\ 0 & -I \end{pmatrix}.$$

This completes the proof.

3.11 Triangular Maps

A *flag* in V is a sequence V_0, \dots, V_r of distinct subspaces such that

$$V_0 \subset V_1 \subset \dots \subset V_r.$$

If $\dim V = n$, then a flag contains at most $n + 1$ subspaces, and a *maximal flag* is a flag with $n + 1$ elements. A maximal flag V_0, \dots, V_n has $V_0 = \{0\}$ and $V_n = V$. If v_1, \dots, v_n is a basis for V and we define $V_0 = \{0\}$ and

$$V_i := \operatorname{span}\{v_1, \dots, v_i\}$$

then V_0, \dots, V_n is a maximal flag. There is a converse to this. Suppose that V_0, \dots, V_n is a maximal flag, and for $i = 1, \dots, n$ choose a non-zero vector w_i in $V_i \setminus V_{i-1}$. Then w_1, \dots, w_n is a basis (as you are invited to prove). Let T be an endomorphism of V . A flag \mathcal{F} is T -invariant if each subspace of \mathcal{F} is T -invariant. If \mathcal{F} is T -invariant, we also say that T fixes \mathcal{F} .

3.11.1 Lemma. *If $\beta = x_1, \dots, x_n$ is a basis for the vector space V and the linear map A fixes the flag associated to β , then the matrix that represents A relative to β is upper triangular.* \square

3.11.2 Theorem. *An endomorphism of a finite-dimensional vector space over an algebraically closed field fixes a maximal flag.*

Proof. We prove the result by induction on $\dim V$. We define a hyperplane in V to be a subspace with dimension $\dim(V) - 1$. It will be enough to prove that any endomorphism of a vector space fixes a hyperplane H for then, by induction, we may assume that $T \upharpoonright H$ fixes a maximal flag of H .

By Lemma 2.5.2, the adjoint T^* of T has an eigenvector in V^* . Choose such an eigenvector f . Then $T^* f = \lambda f$ for some scalar λ , but $T^* f$ is the composition $f \circ T$ and therefore

$$f(Tv) = \lambda f(v),$$

for all elements v of V . This implies that if $f(v) = 0$, then $f(Tv) = 0$ and therefore $\ker f$ is T -invariant.

Since $f \neq 0$, there is a vector v such that $f(v) \neq 0$. If $f(w) \neq 0$ too, then the vector

$$f(w)v - f(v)w$$

lies in $\ker f$, from which it follows that $\ker f$ is a hyperplane. \square

In Section 3.12, we will prove a more concrete version of this result using a variation of the above argument.

- (1) Prove that each maximal flag determines a basis, as described above.
- (2) Prove that if $f \in V^*$, then $\ker f$ is a hyperplane.
- (3) Let S and T be endomorphisms of V that fix the same flag, and suppose $n = \dim V$. Prove that the minimal polynomial of $ST - TS$ divides t^n .

3.12 Triangulations

We prove that if A is a square matrix over \mathbb{C} , then there is a unitary matrix L such that $L^{-1}AL$ is triangular. We have already proved a version of this result for linear mappings (see Section ??) but our argument there did not yield the fact that we could choose L to be unitary.

3.12.1 Theorem. *Let A be an $n \times n$ matrix over \mathbb{C} . Then there is a unitary matrix L such that $L^{-1}AL$ is lower triangular.*

Proof. We proceed by induction on n . Let u_1 be an eigenvector for A^* with eigenvalue θ and let U denote the subspace

$$u_1^\perp = \{x \in \mathbb{C}^n : u_1^* x = 0\}.$$

Then U is A -invariant: if $v \in U$, then

$$u_1^* Av = (A^* u_1)^T v = \theta u_1^* v = 0.$$

Let u_2, \dots, u_n be an orthonormal basis for U . Since $u_1 \notin U$, the vectors u_1, u_2, \dots, u_n form an orthonormal basis for \mathbb{C}^n . If we define the matrix L_1 by

$$L_1 := (u_1 \quad u_2 \quad \cdots \quad u_n)$$

then L_1 is unitary and

$$AL_1 = L_1 \begin{pmatrix} a & 0 \\ b & A_2 \end{pmatrix}.$$

We may assume inductively that there is a unitary matrix M such that $M^{-1}A_2M$ is lower triangular; then $L = L_1M$ is the unitary matrix we need. \square

Suppose that M is an upper triangular $n \times n$ matrix. If $Mv = \theta v$, then $(M - \theta I)v = 0$ and so $M - \theta I$ is not invertible. The matrix $M - \theta I$ is also upper triangular; it is invertible if and only if its diagonal entries are non-zero. We conclude that the eigenvalues of M are precisely the diagonal entries of M . This generalises the fact that the eigenvalues of a diagonal matrix are its diagonal entries.

3.13 The “Fundamental” “Theorem of Algebra”

The fundamental theorem of algebra is the assertion that any polynomial with coefficients from \mathbb{C} has a root in \mathbb{C} . It is equivalent to the claim that every complex matrix has an eigenvector, and we offer a proof of this due to Harm Derksen. The original appears in the American Math. Monthly, and on his web page.

(It has been stated that this result is theorem of analysis, not algebra, and is not fundamental. I tend to agree.)

3.13.1 Theorem. *Every square complex matrix has an eigenvector.*

Before setting out on the proof, some terminology. Let \mathcal{A} be a set with a multiplication defined on it. If $A, B \in \mathcal{A}$, we denote their product by AB . A set \mathcal{A} of endomorphisms of V is an *algebra* if

- (a) \mathcal{A} is a vector space over \mathbb{F} .
- (b) If $A, B \in \mathcal{A}$, then $AB \in \mathcal{A}$.
- (c) There is an element I in \mathcal{A} such that $AI = IA = A$ for all A in \mathcal{A} .

If V is a vector space over \mathbb{F} , then $\text{End}(V)$ is an algebra. If the elements of an algebra \mathcal{A} are endomorphisms of V , it is called an *operator algebra*; if the elements of \mathcal{A} are matrices we call it a *matrix algebra*. The set of all upper triangular matrices is an example of a matrix algebra. The set of strictly upper triangular matrices is not an algebra according to our definition, because it does not contain the identity matrix. An algebra \mathcal{A} is commutative if $AB = BA$ for all A and B in \mathcal{A} . If A is a square matrix, the set of all polynomials in A is a commutative algebra.

We note next that if $f(t)$ is a polynomial over \mathbb{R} with odd degree, then f has a real zero. (This is a comparatively simple exercise in calculus.)

We now start the proof of the theorem. We divide it into a number of lemmas.

3.13.2 Lemma. *If A is an $n \times n$ real matrix and n is odd, then A has an eigenvector.*

Proof. The space \mathbb{R}^n is a direct sum of cyclic subspaces for A . Since n is odd, there is a cyclic subspace U for A with odd dimension d . The minimal polynomial ψ of $A|_U$ has degree d , and therefore there is a real number θ such that $\psi(\theta) = 0$. It follows that A has an eigenvector with eigenvalue θ . \square

3.13.3 Lemma. *If \mathcal{A} is a commutative algebra of real $n \times n$ matrices and n is odd, there is a vector z which is an eigenvector for all matrices in \mathcal{A} .*

Proof. Let A_1, \dots, A_k be a basis for \mathcal{A} . If \mathcal{A} is generated by I , there is nothing to prove, so we may assume $A_1 \neq I$. By the previous lemma, A_1 has an eigenvector z ; let θ be its eigenvalue. The subspaces $\ker(A_1 - \theta I)$ and $\text{im}(A_1 - \theta I)$ are proper non-zero subspaces of \mathbb{R}^n and by the rank theorem,

$$\dim(\ker(A_1 - \theta I)) + \dim(\text{im}(A_1 - \theta I)) = n.$$

Therefore one of these subspaces has odd dimension; we denote it by U .

If $A_1 u = \theta u$, then

$$A_1 A_i u = A_i A_1 u = \theta A_i u$$

and consequently $A_i u \in \ker(A_1 - \theta I)$ if $u \in \ker(A_1 - \theta I)$. If $v = (A_1 - \theta I)w$, then

$$A_i v = A_i(A_1 - \theta I)w = (A_1 - \theta I)A_i w \in \text{im}(A_1 - \theta I).$$

Hence U is invariant under each matrix A_1, \dots, A_k , and so it is invariant under all matrices in \mathcal{A} .

Since U is a proper non-zero subspace of \mathbb{R}^n with odd dimension, it follows by induction that there is a vector in U which is an eigenvector for each matrix in \mathcal{A} . \square

3.13.4 Lemma. *If A is an $n \times n$ complex matrix and n is odd, then A has an eigenvector.*

Proof. Let W denote the vector space of all $n \times n$ Hermitian matrices (which is not an algebra if $n > 1$). We define linear operators L_1 and L_2 by

$$\begin{aligned} L_1(M) &= \frac{1}{2}(AM + MA^*), \\ L_2(M) &= \frac{1}{2i}(AM - MA^*). \end{aligned}$$

If $M = M^*$, then

$$(L_1(M))^* = \frac{1}{2}(AM + MA^*)^* = \frac{1}{2}(MA^* + AM) = L_1(M)$$

and

$$(L_2(M))^* = \frac{1}{-2i}(AM - MA^*)^* = \frac{1}{-2i}(MA^* - AM) = L_2(M).$$

Therefore $L_1, L_2 \in \text{End}(W)$. Also

$$L_1 L_2(M) = \frac{1}{2} \frac{1}{2i} [A(AM - MA^*) + (AM - MA^*)A^*] = \frac{1}{2} \frac{1}{2i} [A^2 M - M(A^*)^2]$$

and

$$L_2 L_1(M) = \frac{1}{2} \frac{1}{2i} [A(AM + MA^*) - (AM + MA^*)A^*] = \frac{1}{2} \frac{1}{2i} [A^2 M - M(A^*)^2].$$

Therefore L_1 and L_2 commute.

Now W is a vector space of dimension n^2 over \mathbb{R} , and n^2 is odd. If we choose a basis for W , the matrices representing L_1 and L_2 relative to this basis have order $n^2 \times n^2$ and they commute. Consequently they have a common eigenvector, and this is an eigenvector for L_1 and L_2 . This eigenvector is a non-zero matrix M such that

$$L_1(M) = \lambda M, \quad L_2(M) = \mu M.$$

Then

$$AM = L_1(M) + iL_2(M) = (\lambda + i\mu)M$$

and this shows that each non-zero column of M is an eigenvector for A . \square

3.13.5 Lemma. *If \mathcal{A} is a commutative algebra of complex $n \times n$ matrices and n is odd, there is a vector z which is an eigenvector for all matrices in \mathcal{A} .*

Proof. We simply apply the proof of Lemma 3.13.3. If A_1, \dots, A_k is a basis for \mathcal{A} and A_1 has an eigenvector, then there is a non-zero proper subspace of \mathbb{C}^n of odd dimension over \mathbb{C} which is invariant under \mathcal{A} . By induction this contains an eigenvector for \mathcal{A} . \square

3.13.6 Lemma. *A square complex matrix has an eigenvector.*

Proof. Assume $n = 2^k n_1$, where n_1 is odd. We prove the lemma by induction on k . Let W denote the space of all matrices M in $\text{Mat}_{n \times n}(\mathbb{C})$ such that $M^T = -M$. We note that

$$\dim(W) = \binom{n}{2}$$

and therefore 2^k does not divide $\dim(W)$. We define two mappings L_1 and L_2 as follows:

$$L_1(M) = AM + MA^T,$$

$$L_2(M) = AMA^T.$$

Then $L_1, L_2 \in \text{End}(W)$ and $L_1 L_2 = L_2 L_1$. Choose a basis for W . The matrices representing L_1 and L_2 relative to this basis commute and have order $\binom{n}{2} \times \binom{n}{2}$.

By induction on k , the algebra generated by these matrices has an eigenvector M ; this is an eigenvector for L_1 and L_2 and we may assume that its eigenvalues are λ and μ respectively. Hence

$$\mu M = AMA^T = A(L_1(M) - AM) = (\lambda A - A^2)M$$

and so

$$(A^2 - \lambda A + \mu I)M = 0.$$

Let z be a non-zero column of M . Then the minimal polynomial of A relative to z is quadratic, and so the A -cyclic subspace generated by z has dimension at most two. Assume that the minimal polynomial ψ of A relative to z is quadratic, and is equal to

$$t^2 - \lambda t - \mu.$$

This quadratic has two roots in \mathbb{C} , and so there are complex numbers θ and τ such that

$$(A - \theta I)(A - \tau I)z = 0.$$

If $(A - \tau I)z = 0$, then z is an eigenvector for A with eigenvalue τ ; if $(A - \tau I)z \neq 0$ then $(A - \tau I)z$ is an eigenvector for A with eigenvalue θ . Thus we have shown that A has an eigenvector. \square

3.14 The Kronecker Product

If A and B are matrices over \mathbb{F} , we construct their *Kronecker product* $A \otimes B$ by replacing the ij -entry of A with

$$A_{i,j}B,$$

for all i and j . We find that

$$(A \otimes B)(u \otimes v) = Au \otimes Bv$$

and, more generally that

$$(A \otimes B)(C \otimes D) = AC \otimes BD,$$

provided only that the products AC and BD are defined. It follows that if x is an eigenvector for A and y is an eigenvector for B , then $x \otimes y$ is an eigenvector for

$A \otimes B$. Consequently the eigenvectors of $A \otimes B$ are just the products $\lambda\mu$, where λ is an eigenvalue of A and μ is an eigenvalue of B . We also have

$$(A \otimes B)^T = A^T \otimes B^T.$$

If X is an $m \times n$ matrix, then $\text{vec}(X)$ is the $mn \times 1$ matrix we get by stacking the columns of X one above the other. In other terms

$$\text{vec}(X) = \sum X_{i,j} e_i \otimes e_j.$$

We have

$$\text{vec}(AX) = (I \otimes A) \text{vec}(X), \quad \text{vec}(XB) = (B^T \otimes I) \text{vec}(X).$$

It follows for example, that there is a matrix X such that

$$AX - XB = C$$

if and only if

$$(I \otimes A - B^T \otimes I) \text{vec}(X) = \text{vec}(C).$$

The eigenvalues of the matrix $I \otimes A - B^T \otimes I$ are the differences $\mu - \lambda$, where λ is an eigenvalue of A and μ is an eigenvalue of B , and therefore it is invertible if and only if A and B have no eigenvalues in common.

Let P be the matrix such that

$$P(x \otimes y) = y \otimes x.$$

Then P maps $U \otimes V$ to $V \otimes U$. If $V = U$, then $P^2 = P$. We say an element u of $V \otimes V$ is *symmetric* if $Pu = u$ and *antisymmetric* if $Pu = -u$. Thus $u \otimes u$ and

$$u \otimes v + v \otimes u$$

are symmetric and

$$u \otimes v - v \otimes u$$

is antisymmetric. (Thus symmetric and antisymmetric elements of $V \otimes V$ are eigenvectors for P , with eigenvalues 1 and -1 respectively.) If A and B belong to $\text{End}(V)$, then

$$P(A \otimes B)P(u \otimes v) = P(A \otimes B)(v \otimes u) = (B \otimes A)(u \otimes v).$$

We also have

$$P \text{vec}(X) = \text{vec}(X^T).$$

- (1) Show that the matrix $P(A \otimes A^T)$ is symmetric.
- (2) Let V be $\text{Mat}_{n \times n}(\mathbb{F})$ and let A be a fixed matrix. If $X \in V$, define the map Ad_A in $\text{End}(V)$ by

$$\text{Ad}_A(X) := AX - XA.$$

If $A^n = 0$, prove that $\text{Ad}_A^{2n} = 0$.

Chapter 4

Orthogonality

We study inner product spaces over \mathbb{R} and \mathbb{C} .

4.1 Properties of Projections

Let U be subspace of the inner product space V . Then orthogonal projection onto U is a function P from V to itself such that, for all v in V , we have $v - P(v) \in U^\perp$. We establish a number of properties of P , the most important of which is that it is linear.

4.1.1 Lemma. *Let P be the orthogonal projection of V onto U . Then P is linear mapping and*

(a) $\text{im}(P) = U$.

(b) $\text{ker}(P) = U^\perp$.

(c) $P^2 = P$.

(d) If $v, w \in V$, then $\langle v, Pw \rangle = \langle Pv, w \rangle$.

Proof. Suppose $v, w \in V$. Then $v - P(v)$ and $w - P(w)$ both belong to U^\perp , whence

$$(v + w) - (P(v) + P(w)) = (v - P(v)) + (w - P(w)) \in U^\perp.$$

Since $P(v) + P(w) \in U$, this implies that $P(v) + P(w)$ is the orthogonal projection of $v + w$ onto U . Therefore P is linear.

Since $Pv \in U$ for all v in V , we see that $\text{im}(P) \subseteq U$, and since $Pu = u$ for all u in U , we see that $\text{im}(P) = U$ and $P^2 = P$. This proves (a) and (c).

For (b), we note that if $P(v) = 0$ then $v \in U^\perp$. On the other hand, $v - P(v) \in U^\perp$ and so if $v \in U^\perp$ then $P(v) \in U^\perp$. Since $P(v) \in U$, this implies that $P(v) = 0$.

Finally, for any vectors v and w we have

$$\langle v - Pv, Pw \rangle = 0, \quad \langle Pv, Pw - w \rangle = 0.$$

Summing these two expressions yields

$$0 = \langle v, Pw \rangle - \langle P, v \rangle Pw + \langle Pv, Pw \rangle - \langle Pv, w \rangle,$$

whence (d) follows. \square

Linear mappings P such that $P^2 = P$ are called *idempotent*. If $\langle v, Pw \rangle = \langle Pv, w \rangle$ for all v and w , we say P is *self adjoint* with respect to the inner product.

4.2 Matrices Representing Projections

If we are working in Euclidean space— \mathbb{R}^n with dot product—then we can give an explicit formula for the matrix representing orthogonal projection.

4.2.1 Lemma. *Let V be \mathbb{R}^n equipped with dot product, and let U be a subspace of V with dimension k . If B is an $n \times k$ matrix whose columns form a basis for U , the matrix representing orthogonal projection on U is*

$$B(B^T B)^{-1} B^T.$$

Proof. We offer two proofs. The first is a simple verification that the quoted formula is correct. First you may easily verify that $B(B^T B)^{-1} B^T$ is symmetric. Then we compute that

$$(I - B(B^T B)^{-1} B)B = B - B(B^T B)^{-1} B^T B = B - B = 0$$

and therefore

$$\begin{aligned} (v - B(B^T B)^{-1} B^T v)^T B &= v^T (I - B(B^T B)^{-1} B)^T B \\ &= v^T (I - B(B^T B)^{-1} B)B \\ &= 0. \end{aligned}$$

So, if $u := B(B^T B)^{-1} B^T v$, then $v - u$ is orthogonal to each column of B . Hence it lies in U^\perp , and therefore u is the orthogonal projection of v onto U .

A difficulty with the previous argument is that it gives no indication how we found the matrix $B(B^T B)^{-1} B^T$ in the first place. We outline the reasoning. Suppose Q is the matrix representing orthogonal projection on U . Then $\text{rk } Q = k$ and by Theorem ?? we can write

$$Q = AB^T,$$

where A and B are $n \times k$ matrices with rank k . Since $\text{col}(A) = \text{col}(Q) = U$, the columns of A form a basis for U . Since the columns of A are linearly independent, if $Ax = 0$ then $x = 0$. Therefore $\ker B^T = \ker P = U^\perp$ and consequently $\text{col}(B) \subseteq U^{\perp\perp} = U$. As $\text{rk } B = k$, this shows that $\text{col } B = U$. Since each column of A lies in U , we have

$$A = QA = AB^T A$$

and therefore $B^T A = I$. On the other hand, the columns of B form a basis for U , so each column of A is a linear combination of columns of B , and therefore there is a $k \times k$ matrix M such that $A = BM$. If $B^T A = I$, this implies that $B^T BM = I$ and so $M = (B^T B)^{-1}$. Accordingly

$$Q = B(B^T B)^{-1} B^T. \quad \square$$

(It might be a useful exercise to identify where in this argument we have used that our inner product is the dot product.)

To sum up, we have two ways to compute the orthogonal projection of a vector v onto a subspace U . If we are given an orthogonal basis for U , we can use (??). If we are working in \mathbb{R}^n with dot product and given a basis for U , we could construct $Q = B(B^T B)^{-1} B^T$, in which case the answer is Qv .

- (1) If $\langle \cdot, \cdot \rangle$ is the dot product, show that (??) implies that $P = P^T$.
- (2) Suppose B and C are $n \times k$ matrices with rank k and the same column space. Prove that $B(B^T B)^{-1} B = C(C^T C)^{-1} C$.
- (3) Let u_1, \dots, u_k be an orthogonal basis for the subspace U . Show that the matrix representing orthogonal projection on U is equal to

$$\sum_{i=1}^k \langle u_i, u_i \rangle^{-1} u_i u_i^T.$$

4.3 Least Squares

We consider a version of the least squares problem. Let W be an $m \times n$ matrix where $m < n$, and $\text{rk}(W) = m$. Then the system of equations

$$Wx = v \tag{4.3.1}$$

will have infinitely many solutions, but some may suit us better than others. For example, in the control theory setting of Chapter 11, a solution x to an equation of the form $Wx = v$ represented a sequence of inputs that would drive our system to a chosen state. In this case, $x^T x$ corresponds to the power that this sequence would require, and it would be very natural to seek to minimize it. Thus we want to find the solution to (4.3.1) with minimum squared length.

Suppose that x is any solution to (4.3.1), and let \bar{x} be the orthogonal projection of x on $\text{col}(W^T)$. Then $x - \bar{x}$ is orthogonal to $\text{col}(W^T)$, and therefore $W(x - \bar{x}) = 0$. Hence

$$W\bar{x} = Wx = v.$$

If y is another solution to (4.3.1), then $Wy = W\bar{x}$ and so $y - \bar{x}$ is in the null space of W . Consequently $y - \bar{x}$ is orthogonal to \bar{x} and

$$\|y\|^2 = \|y - \bar{x}\|^2 + \|\bar{x}\|^2.$$

Thus \bar{x} is the solution to (4.3.1) with minimum norm.

How can we compute \bar{x} ? If we can assume that the rows of W are linearly independent then, by Lemma 4.2.1, the matrix representing orthogonal projection onto $\text{col } W^T$ is

$$Q := W^T(WW^T)^{-1}W$$

and our solution is Qx . However we do not need to find x ; we have

$$Qx = W^T(WW^T)^{-1}Wx = W^T(WW^T)^{-1}v,$$

and we can proceed as follows: given v , solve the system

$$WW^T z = v,$$

the desired solution is $W^T z$. (This approach avoids the need to compute the inverse of WW^T . Computing an inverse explicitly is rarely worth the trouble. It may also pay to avoid computing WW^T , but we digress....)

In ??, we will develop a general method for least squares problems, which does not require that the rows of W are linearly independent.

4.4 Orthogonal Polynomials

Let V be the space of all real polynomials, or the vector space of polynomials with degree at most n . Assume V is equipped with an inner product such that

$$\langle p, xq \rangle = \langle xp, q \rangle,$$

and, if $p(x)$ is non-negative and not zero, then $\langle 1, p \rangle > 0$. All our examples have these properties.

If we apply Gram-Schmidt to the basis of V formed by the powers of x , we obtain a sequence of polynomials $(p_r)_{r \geq 0}$, where p_r has degree r . A *sequence of orthogonal polynomials* is an orthogonal set of polynomials $(p_r)_{r \geq 0}$, where p_r has degree r (and $p_0 \neq 0$). If we multiply each member of a sequence of orthogonal polynomials by a non-zero scalar, the result is still a sequence of orthogonal polynomials.

4.4.1 Lemma. *The sequence of polynomials $(p_r)_{r \geq 0}$ is an orthogonal basis if and only if p_r is non-zero and is orthogonal to all polynomials of degree less than r .* \square

4.4.2 Lemma. *Let $(p_r)_{r \geq 0}$ be a sequence of orthogonal polynomials. If $p_r(x) = a(x)b(x)$, where a and b are polynomials and $b(x) \geq 0$ for all x , then b is constant.*

Proof. We have

$$\langle p_r, a \rangle = \langle ab, a \rangle = \langle 1, a^2 b \rangle.$$

Now $a^2 b$ is non-zero and non-negative and therefore $\langle 1, a^2 b \rangle > 0$. But, if the degree of b is positive, then the degree of a is less than r and, by the previous lemma, $\langle p_r, a \rangle = 0$. We conclude that b must be constant. \square

4.4.3 Theorem. *If p is a member of a sequence of orthogonal polynomials, its zeros are real and simple.*

Proof. Suppose θ is a complex zero of p . Then its complex conjugate $\bar{\theta}$ is also a zero of p and therefore the real quadratic polynomial

$$(x - \theta)(x - \bar{\theta})$$

divides p . Since this quadratic has two complex roots and is monic, it is non-negative. By the previous lemma, it cannot divide p . This proves the first claim.

For the second, note that $(x - \theta)^2$ is non-negative and the same technique yields that this cannot divide p . \square

- (1) Let $\langle p, q \rangle := \int_0^1 p(x)q(x) dx$. Show that if p is a member of the sequence of orthogonal polynomials associated to this inner product, all zeros of p lie in the interval $(0, 1)$.
- (2) Suppose p_r and p_{r+1} are consecutive members of a sequence of orthogonal polynomials. Show that they cannot have a common zero.

4.5 The Three-Term Recurrence

We provide an easier way to construct families of orthogonal polynomials. The key is to note that

$$\langle xp_r, p_j \rangle = 0$$

if $j \notin \{r-1, r, r+1\}$. For if $j < r-1$ then xp_j has degree less than r , and therefore

$$\langle xp_r, p_j \rangle = \langle xp_r, p_j \rangle = 0.$$

If $j > r+1$ then similarly p_j is orthogonal to xp_r .

4.5.1 Theorem. *Let $(p_r)_{r \geq 0}$ be a sequence of monic orthogonal polynomials. Then*

$$p_{n+1} = (x - a_n)p_n - b_n p_{n-1},$$

where $a_n = \langle xp_n, p_n \rangle / \langle p_n, p_n \rangle$ and $b_n = \langle p_n, p_n \rangle / \langle p_{n-1}, p_{n-1} \rangle$.

Proof. From our remarks just above, xp_n is a linear combination of p_{n-1} , p_n and p_{n+1} . Thus we may write

$$xp_n = \gamma p_{n+1} + \alpha p_n + \beta p_{n-1}.$$

Here

$$\gamma = \frac{\langle xp_n, p_{n+1} \rangle}{\langle p_{n+1}, p_{n+1} \rangle}.$$

Since p_{n+1} is monic, $xp_n = p_{n+1} - q$, where q has degree less than n . So

$$\langle xp_n, p_{n+1} \rangle = \langle p_{n+1}, p_{n+1} \rangle$$

and therefore $\gamma = 1$.

Next we see that $\alpha = \langle xp_n, p_n \rangle / \langle p_n, p_n \rangle$ and

$$\beta = \frac{\langle xp_n, p_{n-1} \rangle}{\langle p_{n-1}, p_{n-1} \rangle}.$$

Arguing as before,

$$\langle xp_n, p_{n-1} \rangle = \langle p_n, p_n \rangle$$

and this leads to the stated expression for b_n . □

One consequence of the formulas for the coefficients in this recurrence is that $b_n > 0$ for all n .

There is another way of stating the last result. Let $(p_r)_{r \geq 0}$ be a monic sequence of orthogonal polynomials. Let M_x denote the linear transformation that maps a polynomial p to xp . Then the matrix representing M_x with respect to the basis $(p_r)_{r \geq 0}$ is

$$\begin{pmatrix} a_0 & b_1 & & & \\ 1 & a_1 & b_2 & & \\ & 1 & a_2 & b_3 & \\ & & \ddots & \ddots & \ddots \end{pmatrix}$$

This is an example of a *tridiagonal matrix*.

4.6 Numerical Integration

We want to compute definite integrals of the form

$$\int_a^b f(t) w(t) dt.$$

Here $w(t)$ is a *weight function*. For example if the interval of integration is $[0, \infty]$, then we may use $w(t) = e^{-t}$. But for now we take $w(t)$ to be identically 1, and the interval of integration will be $[0, 1]$. So all we want is

$$\int_0^1 f(t) dt.$$

The problem is that we do not know the anti-derivative of f , and so we seek a procedure that will produce a reasonably accurate answer in reasonable time.

There are many possible notions of what ‘reasonably accurate’ might mean. Before we discuss this, we specify the sort of procedure we want in more detail. The first thing is to note that the map

$$\mathcal{S} : f \mapsto \int_0^1 f(t) dt.$$

is a linear map from the space $C[0, 1]$ of continuous functions on $[0, 1]$ to \mathbb{R} . Hence it is an element of the dual space $C[0, 1]^*$. This has the property that if $f \geq 0$, then $\mathcal{S}(f) \geq 0$ —it maps non-negative functions to non-negative real

numbers. In this context, the elements of $C[0, 1]^*$ are known as *linear functionals*, and we say a linear functional is *non-negative* if it maps non-negative functions to non-negative numbers.

There are many other non-negative linear functionals, and amongst the simplest are the *evaluation maps* e_a , for $a \in \mathbb{R}$, given by

$$e_a(f) := f(a).$$

Our aim is to choose an increasing sequence of *nodes* $\theta_1, \dots, \theta_n$ and a sequence of *weights* w_1, \dots, w_n , such that the linear functional

$$\mathcal{Q} := \sum_i w_i e_{\theta_i}$$

is a good approximation to \mathcal{S} . We call a linear functional of this form a *quadrature scheme*. Define the *degree of precision* of \mathcal{Q} to be the greatest integer k such that

$$\mathcal{Q}(p) = \int_0^1 p(t) dt$$

for all polynomials p with degree at most k .

By way of example, if \mathcal{Q} has degree of precision 1, then

$$\mathcal{Q}(1) = 1 \quad \mathcal{Q}(t) = \frac{1}{2}.$$

These hold if and only if

$$\sum_i w_i = 1, \quad \sum_i w_i \theta_i = \frac{1}{2}.$$

It is easy to find nodes and weights for which these conditions hold.

We will be more greedy. Suppose we are given nodes a_1, \dots, a_n , and that we try to find weights to go with them. Let p_1, \dots, p_n be the Lagrange interpolating polynomials at the given nodes. Thus

$$p_i(a_j) = \delta_{i,j}.$$

Then

$$\mathcal{Q}(p_i) = \sum_j w_j p_i(a_j) = w_i$$

and, if we want degree of precision at least $n - 1$, we will need

$$w_i = \int_0^1 p_i(t) dt.$$

There is one problem: there are good reasons to require that the weights w_i be non-negative, and it is not clear how to choose the nodes to ensure this in general.

We can go further if we use orthogonal polynomials. Define an inner product on $C[0, 1]$ by

$$\langle p, q \rangle = \int_0^1 p(t)q(t) dt.$$

Let p_0, \dots, p_n be the first $n + 1$ members of the corresponding family of orthogonal polynomials, and let $\theta_1, \dots, \theta_n$ be the zeroes of p_n in increasing order. (We know by Theorem 4.4.3 that these zeroes are real and distinct.) Using the Lagrange interpolating polynomials, we compute the weights w_i for a quadrature scheme with degree of precision at least $n - 1$. Then, as Gauss first noted, a miracle occurs: the degree of precision of our scheme is $2n - 1$.

We verify this. Suppose f is a polynomial with degree at most $2n - 1$. By the Euclidean algorithm, there are polynomials q and r , both with degree at most $n - 1$, such that

$$f(t) = q(t)p_n(t) + r(t).$$

Now

$$\int_0^1 f(t) dt = \langle 1, f(t) \rangle = \langle 1, q(t)p_n(t) \rangle + \langle 1, r(t) \rangle.$$

Since q has degree less than n ,

$$\langle 1, q(t)p_n(t) \rangle = \langle q(t), p_n(t) \rangle = 0,$$

and therefore

$$\int_0^1 f(t) dt = \int_0^1 r(t) dt.$$

Because the degree of $r(t)$ is at most $n - 1$, the integral on the right can be computed exactly using our (well, Gauss's) quadrature scheme. Hence this scheme has degree of precision at least $2n - 1$. However it is exactly $2n - 1$, because

$$\mathcal{Q}(p_n(t)^2) = 0 < \int_0^1 p_n(t)^2 dt.$$

Chapter 5

Eigenthings

In this chapter we undertake a study of questions related to existence of eigenvectors and eigenvalues. Our focus is on self-adjoint operators, because that is where eigenvalues are most useful.

5.1 Self-Adjoint Operators

If S is an operator on an inner product space V , we define the *adjoint* S^* of S to be an operator such that, for all u, v in V we have

$$\langle S^* u, v \rangle = \langle u, Sv \rangle.$$

It is an easy exercise to show that, if it exists, the adjoint is unique. If V is \mathbb{C}^n and the inner product is the usual complex dot product, and M is a matrix representing S on V , then S^* is represented by the conjugate transpose of M , which we usually denote by M^* .

By way of a second example, if V is the vector space of all real polynomials and

$$\langle p, q \rangle := \int_a^b p(x)q(x)w(x)dx$$

then

$$\langle xp, q \rangle = \langle p, xq \rangle.$$

Hence the operation of multiplication by x is a self-adjoint linear mapping of V . (This is why the theory of orthogonal polynomials is so rich.)

We turn to the existence question. There is a notational difficulty that arises because, outside the context of inner product space, the dual of a linear mapping is often referred as the adjoint (with good reason). We will temporarily use S^d to denote the dual of a linear mapping.

There are two steps to the existence proof. Assume V is an inner product space. If $a \in V$ then we have a map θ_a in V^* given by

$$\theta_a(v) = \langle a, v \rangle.$$

This is linear, and is an isomorphism from V to its dual V^* . If $S \in \text{End } V$ then, by the definition of the dual,

$$\langle u, Sv \rangle = (\theta_a \circ S)(b) = (S^d \theta_a)(b).$$

Now $S^d \circ \theta_a \in V^*$ and, since the map $a \mapsto \theta_a$ is an isomorphism, there is a vector $S^*(a)$ in V such that

$$S^d \circ \theta_a = \theta_{S^*(a)}.$$

As the notation suggests, and as you should prove, the map $a \mapsto S^*(a)$ is linear. We say that S^* is the *adjoint* of S .

5.2 Diagonalizability

We prove that self-adjoint operators are diagonalizable, and more.

5.2.1 Theorem. *Let S be a self-adjoint operator on the inner product space V . Then*

- (a) *The minimal polynomial of S has only simple zeros.*
- (b) *Eigenvectors of S with distinct eigenvalues are orthogonal.*
- (c) *The eigenvalues of S are real.*
- (d) *S is diagonalizable.*

Proof. Assume by way of contradiction that the minimal polynomial $\psi(t)$ of S has a multiple root. Then there is a proper divisor $p(t)$ of $\psi(t)$ such that $\psi(t)$ divides $p(t)^2$. Suppose $S_1 := p(S)$. Then $S_1 \neq 0$ because p is a proper divisor of ψ , but $S_1^2 = 0$. If $v \in V$, then $S_1^2 v = 0$ and

$$0 = \langle v, S_1^2 v \rangle = \langle S_1 v, S_1 v \rangle,$$

whence we see that $S_1 v = 0$. Since v can be chosen arbitrarily in V , this implies that $S_1 = 0$, our contradiction.

Next, suppose

$$Su = \theta u, \quad Sv = \tau v.$$

Then

$$\tau \langle u, v \rangle = \langle u, \tau v \rangle = \langle u, Sv \rangle = \langle Su, v \rangle = \langle \theta u, v \rangle = \theta \langle u, v \rangle;$$

and we conclude that either $\theta = \tau$, or $\langle u, v \rangle = 0$.

It follows from (a) that S is diagonalizable, and therefore V is a direct sum of S -invariant subspaces V_i on each of which S acts as multiplication by a scalar θ_i .

On the other hand, the proof of Corollary 5.4.2 shows that each S^*S -invariant subspace of V contains an eigenvector for S^*S with a non-negative real eigenvalue. Since $S^*S = S^2$, any S -invariant subspace is S^*S -invariant. Suppose the eigenvalue of S^2 on V_i is σ_i . Then $\sigma_i = \theta_i^2$, and therefore θ is real. \square

5.2.2 Corollary. *Suppose S is a self-adjoint operator on the inner product space V . Then there is an orthogonal basis for V formed of eigenvectors for S .* \square

5.3 Diagonalizability, Again

We offer a second proof that for any self-adjoint operator on a finite dimensional space, there is an orthogonal basis for the space that consists of eigenvectors. They is the following result.

5.3.1 Lemma. *Let S be a self-adjoint operator on the inner product space V . If U is an S -invariant subspace of V , then U^\perp is S -invariant.*

Proof. If $v \in U^\perp$, then

$$\langle Sv, u \rangle = \langle v, Su \rangle$$

and therefore Sv lies in U^\perp . \square

This makes everything easy. Suppose S is self-adjoint and λ is a zero of its minimal polynomial. Then there is an eigenvector z associated with λ . Its span U is certainly S -invariant, and hence U^\perp is S -invariant. The restriction of S to U^\perp is self-adjoint (prove it) and therefore by induction on the dimension, U^\perp has an orthogonal basis formed from eigenvectors. This basis together with z provides the basis of V that we need. \square

5.4 Eigenvectors and Optimization

We present a result which may appear to be of limited interest, but it provides an important reason why we should be interested in eigenvectors. It also illustrates how self-adjoint operators can arise in practice.

5.4.1 Lemma. *Let L be a linear map from \mathbb{R}^n to \mathbb{R}^m , let U be a subspace of \mathbb{R}^n , and let u be a unit vector in U such that $\|Lu\|$ is maximal. If $h \in U$ and $h^T u = 0$, then $h^T L^T Lu = 0$.*

Proof. We have

$$\|L(u + th)\|^2 = (u + th)^T L^T L(u + th) = u^T L^T Lu + 2tu^T L^T Lh + t^2 h^T L^T Lh.$$

Since $\langle u, h \rangle = 0$, we have

$$\|u + th\|^2 = \|u\|^2 + t^2 \|h\|^2 = 1 + t^2 \|h\|^2.$$

Assuming that t is small enough that t^2 is negligible, we find that

$$\|Lu\|^2 - \frac{\|L(u + th)\|^2}{\|u + th\|^2} \approx -2th^T L^T Lu.$$

We may choose t to be positive or negative; as we have chosen the unit vector u in U to maximize $\|Lu\|$ it follows that if h is orthogonal to u , then $h^T L^T Lu = 0$, and therefore Lu and h are orthogonal. \square

Now we present the application of this lemma to eigenvectors.

5.4.2 Corollary. *Let L be a linear map from \mathbb{R}^n to \mathbb{R}^m , let U be a subspace of \mathbb{R}^n , and let u be a unit vector in U such that $\|Lu\|$ is maximal. If U is $L^T L$ -invariant, then u is an eigenvector of $L^T L$, and its eigenvalue is non-negative and real.*

Proof. Suppose u is as stated. From the previous lemma we see that if $h \in U$ and $h \in u^\perp$, then $h^T L^T Lu = 0$. Therefore

$$U \cap u^\perp \subseteq (L^T L)^\perp,$$

from which we have

$$L^T Lu \in \text{span}(u) + U^\perp.$$

Therefore $L^T Lu = \theta u + v$, where $v \in U^\perp$. But U is $L^T L$ -invariant, and therefore $L^T Lu \in U$. Hence

$$L^T Lu - \theta u \in U^\perp \cap U = \{0\}$$

and so u is an eigenvector for $L^T L$. \square

Obviously \mathbb{R}^n itself is $L^T L$ -invariant, and thus it follows that if u is a unit vector in \mathbb{R}^n that maximizes $u^T L^T L u$, then u is an eigenvector for $L^T L$. Since the associated eigenvalue is the maximum value of a non-negative real function, the final claim holds. \square

We consider one important case where we are interested in maximizing $\|Lu\|$ over unit vectors. Let A be an $n \times n$ invertible matrix and consider the system of linear equations

$$Ax = b. \quad (5.4.1)$$

If z is a vector then the solution to $Ax = b + z$ is $A^{-1}b + A^{-1}z$. Thus we may say that an error z in b leads to an error $A^{-1}z$ in the solution to (5.4.1).

Which vector z leads to the greatest error? It is clear that if, for example, we replace z by $2z$ then the error is doubled, thus it makes sense to consider

$$\max_{\|z\|=1} \|A^{-1}z\|.$$

From our considerations above, the maximum value of this occurs when z is an eigenvector of

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

The magnitude of the error will be given by the eigenvalue associated with z . We will see that the eigenvalues of AA^T are real and positive. If the matrix M is invertible, then θ is an eigenvalue of M if and only if θ^{-1} is an eigenvalue of M^{-1} . It follows that the solution of (5.4.1) will be most sensitive to errors in b when the least eigenvalue of AA^T is small.

5.5 The Singular Value Decomposition

If the $m \times n$ matrix A has rank k , then it can be shown that there is an $m \times k$ matrix X and a $k \times n$ matrix Y such that $\text{rk}(X) = \text{rk}(Y) = k$ and $A = XY^T$. When we work over \mathbb{R} (or \mathbb{C}), we can prove a somewhat stronger version of this, known as the *singular value decomposition*. This is extremely important in practice.

5.5.1 Theorem. *Let A be a non-zero real matrix with rank k . Then $A = Y\Sigma X^T$, where*

(a) $X^T X = I_k$,

(b) Σ is a $k \times k$ diagonal matrix Σ with positive diagonal entries,

(c) $Y^T Y = I_k$.

Proof. Assume A is $m \times n$. Using induction on k , we construct an orthonormal subset x_1, \dots, x_k of \mathbb{R}^n and an orthonormal subset y_1, \dots, y_k of \mathbb{R}^m such that $y_i = \sigma_i A x_i$ and

$$A = \sum_{i=1}^k \sigma_i y_i x_i^T.$$

This is equivalent to the statement of the theorem.

Let U_0 denote \mathbb{R}^n and let x_1 be a unit vector in U_0 such that $\|A x_1\|$ is maximal. Set σ_1 equal to $\|A x_1\|$ and define

$$x_1 := x, \quad y_1 = \sigma_1^{-1} A x_1.$$

Let U_1 denote x_1^\perp . By Lemma 5.4.1 we see that if $h^T x_1 = 0$, then $h^T A^T A x_1 = 0$ and consequently $A(U_1^\perp) \subseteq A(U_1)^\perp$.

Suppose

$$A_1 := A - \sigma_1 y x^T.$$

Since y lies in the column space of A , we see that $\text{col}(A_1) \subseteq \text{col}(A)$. Since $A \neq 0$ we see that $x_1 \neq 0$ and $y_1 \neq 0$. Therefore $A x_1 \neq 0$, but

$$A_1 x_1 = A x_1 - \sigma_1 y x_1^T x_1 = A x_1 - \sigma_1 y = 0.$$

Consequently $\text{rk}(A_1) < \text{rk}(A)$. As $\text{rk}(\sigma_1 y x_1^T) = 1$ it follows that $\text{rk}(A_1) = k - 1$.

Note next that if $x \in U_1$, then $A x = A_1 x$ and so A and A_1 agree on U_1 . Working now with A_1 and U_1 , we conclude by induction on k that there are orthogonal unit vectors x_2, \dots, x_k in \mathbb{R}^n and orthogonal unit vectors y_2, \dots, y_k in \mathbb{R}^m , such that $y_i = A_1 x_i$ and, if $\sigma_i := \|A_1 x_i\|$, then

$$A_1 = \sum_{i=2}^k \sigma_i y_i x_i^T.$$

Our theorem follows immediately. \square

In numerical work, the following alternative version of the singular value decomposition may be more useful. (It does not assume we know the rank of A .)

5.5.2 Corollary. *If A is a square matrix, there are orthogonal matrices X and Y , and a non-negative diagonal matrix Σ such that $A = Y \Sigma X^T$.* \square

The matrices Y and X in the singular value decomposition $Y\Sigma X^T$ of A are not unique in any useful sense. However Σ is determined up to a permutation. Its entries are known as the *singular values* of A ; there are usually denoted by $\sigma_1, \dots, \sigma_n$, with the assumption that they form a non-increasing sequence.

The easiest way to see that the singular values are determined by A is to verify that their squares are the eigenvalues of AA^T . To show this, note that

$$AA^T = Y\Sigma X^T X\Sigma Y^T = Y\Sigma^2 Y^T,$$

and therefore

$$AA^T Y = Y\Sigma^2.$$

It follows from this that the columns of Y are eigenvectors for AA^T , and the diagonal entries of Σ^2 are its eigenvalues.

In a similar fashion we can show that the squares of the singular values of A are the eigenvalues of $A^T A$. Hence we see that AA^T and $A^T A$ have the same eigenvalues. (This actually holds over any field, although the proof at hand only works over \mathbb{R} or, with modest extra effort, over \mathbb{C} .)

- (1) Prove Corollary 5.5.2.
- (2) Compute the singular values of a companion matrix. (You may work with either CC^T or $C^T C$, but one is significantly easier. First show that all but two of the singular values are equal to 1.)
- (3) Show that the sum of the singular values of a square matrix is a norm.
- (4) If $\sigma_1(A)$ denote the largest singular value of A , show that it is a norm.

.

5.6 Least Squares

We consider the system of linear equations

$$Ax = b \tag{5.6.1}$$

where A is $m \times n$. In ?? we considered the case where the rows of A are linearly independent. Then the columns of A span \mathbb{R}^m , and we want the vector x with

minimum norm such that $Ax = b$. The second, and more commonly met situation, is when the columns of A are linearly independent, and we want the vector x such that $\|b - Ax\|^2$ is minimal.

We draw attention to one difficulty. It is in fact a non-trivial numerical problem to determine the rank of a real matrix, and so it may not be easy to verify that the rows or columns of A are linearly independent. In fact, the best way to determine the rank in finite precision arithmetic is to use the singular value decomposition $A = Y\Sigma X^T$, since $\text{rk}(A) = \text{rk}(\Sigma)$. (Thus determining the rank of A is reduced to determining the rank of a diagonal matrix; in the presence of rounding errors and uncertainties in the data, this still may require thought.) But rather than using the singular value decomposition just to get the rank of A , we can use it to solve the least squares problem.

5.6.1 Lemma. *Let A be an $m \times n$ real matrix with singular value decomposition $A = Y\Sigma X^T$, where Σ is $k \times k$ and invertible. Then the vector z of minimum norm, such that $b - Az$ has minimum norm is given by*

$$z = X\Sigma^{-1}Y^T b.$$

Proof. We note that the columns of Y form an orthonormal basis for $\text{col}(A)$, whence the matrix representing projection onto $\text{col}(A)$ is YY^T . Similarly, the columns of X form an orthonormal basis for $\text{col}(A^T)$, and therefore XX^T is the matrix representing projection onto $\text{col}(A^T)$.

Consequently $y = YY^T b$ is the vector in $\text{col}(A)$ closest to b . Suppose $Ax = y$. Then

$$Y\Sigma X^T x = YY^T b$$

and, multiplying both sides on the left by Y^T , we have

$$X^T x = \Sigma^{-1}Y^T b.$$

Now $XX^T x$ is the projection of x onto $\text{col}(A^T)$, and accordingly

$$z = X\Sigma^{-1}Y^T b$$

is the vector of minimum norm such that Az is closest to b . □

5.7 Legendre Polynomials

Let V be $\text{Pol}(\mathbb{R})$, the vector space of all real polynomials, with inner product

$$\langle p, q \rangle = \int_{-1}^1 p(t)q(t) dt.$$

Define a linear mapping $L: V \rightarrow V$ by

$$L(p) = (1 - t^2)p'' - 2tp'.$$

If $n \geq 2$ then

$$L(t^n) = (1 - t^2)n(n-1)t^{n-2} - 2nt^n = -n(n+1)t^n + n(n-1)t^{n-2}. \quad (5.7.1)$$

It follows that

$$\langle t^m, Lt^n \rangle = \int_{-1}^1 (n(n-1)t^{m+n-2} - n(n+1)t^{m+n}) dt;$$

when $m+n$ is odd the integral here is zero, if $m+n$ is even then it is

$$\left[\frac{2n(n-1)}{m+n-1} - \frac{2n(n+1)}{m+n+1} \right] = -\frac{4mn}{(m+n)^2 - 1}.$$

Hence

$$\langle t^m, Lt^n \rangle = \langle Lt^m, t^n \rangle$$

for all m and n . It follows that for any polynomials p and q ,

$$\langle p, Lq \rangle = \langle Lp, q \rangle,$$

and therefore L is self-adjoint. (This can also be proved directly using integration by parts.)

It follows that the eigenvalues of L are real, and eigenvectors with distinct eigenvalues are orthogonal with respect to the above inner product. It is not hard to determine the eigenvalues of L . From (5.7.1) we see that $\text{Pol}_n(\mathbb{R})$ is L -invariant and further, if L_n denotes the restriction of L to $\text{Pol}_n(\mathbb{R})$ and $\beta = \{1, t, \dots, t^n\}$ is the standard basis for $\text{Pol}_n(\mathbb{R})$, then

$$[L_n]_\beta = \begin{pmatrix} 0 & 0 & 2 & & \\ & -2 & 0 & 6 & \\ & & -6 & 0 & 12 \\ & & & -12 & 0 \\ & & & & \ddots \end{pmatrix}$$

This is a triangular matrix, and reveals that the eigenvalues of L_n are the integers $-m(m-1)$ for $m = 1, \dots, n$.

As the eigenvalues are distinct, each eigenspace is 1-dimensional and is thus spanned by a polynomial. The polynomial with eigenvalue $-m(m-1)$ will have degree m and is a solution of *Legendre's equation*:

$$(1-t^2)p'' - 2tp' + m(m-1)p = 0.$$

We call p_m the *Legendre polynomial* of degree m . The first five Legendre polynomials are as follows:

$$\begin{aligned} p_0 &= 1 \\ p_1 &= t \\ p_2 &= 3t^2 - 1 \\ p_3 &= 5t^3 - 3t \\ p_4 &= 35t^4 - 30t^2 + 3. \end{aligned}$$

It makes no harm if we replace p_i by any non-zero scalar multiple of itself, and it is customary to choose the multiple so that $p_i(1) = 1$. (But we have not done that here.)

There are a number of related examples (of self-adjoint linear operators on $P(\mathbb{R})$). We summarize some of them here. The numbers $\lambda_0, \lambda_1, \dots$ are the eigenvalues of the operator.

(a) **Chebyshev.**

$$Lp = (1-t^2)p'' - tp'; \quad \langle p, q \rangle = \int_{-1}^1 p(t)q(t) \frac{dt}{\sqrt{1-t^2}}; \quad \lambda_n = -n^2.$$

(a) **Laguerre.**

$$Lp = tp'' + (1-t)p'; \quad \langle p, q \rangle = \int_0^\infty p(t)q(t)e^{-t} dt; \quad \lambda_n = -n.$$

(a) **Hermite.**

$$Lp = p'' - tp'; \quad \langle p, q \rangle = \int_{-\infty}^\infty p(t)q(t)e^{-t^2/2} dt; \quad \lambda_n = -n.$$

In general, if

$$Lp = fp'' + gp'$$

then we may write

$$Lp = w^{-1}(wf p)'$$

where

$$w(t) = \frac{1}{f(t)} \exp \int_{\alpha}^t \frac{g(u)}{f(u)} du.$$

(The value of the lower limit α in this integral will be determined by context.)

Then L is self-adjoint relative to the inner product

$$\langle p, q \rangle = \int p(t)q(t) w(t) dt.$$

To see this, compute in outline as follows:

$$\begin{aligned} \langle Lp, q \rangle &= \int q(wf p)'\ dt = - \int wf p' q' dt = - \int wf q' p' dt \\ &= \int p(wf q)'\ dt \\ &= \langle p, Lq \rangle. \end{aligned}$$

For this computation to be accurate, $f(t)w(t)$ must vanish at the endpoints of the interval over which we integrate.

The eigenvectors of L will be polynomials only if $w(t)$ satisfies further restrictions.

5.8 Computing Eigenvalues

How do people really compute the eigenvalues of symmetric matrices? They do not use the method offered in most introductory linear algebra course—compute the characteristic polynomial, find its zeros—that is probably the fourth best method. Here we outline the second best.

So, suppose A is a real symmetric $n \times n$ matrix. We want to find an orthogonal matrix L such that $L^T A L$ is diagonal. What we will actually do is to describe how to find a sequence of orthogonal matrices S_1, \dots, S_r such that all off-diagonal entries of

$$S_r^T \cdots S_1^T A S_1 \cdots S_r$$

are very small, we can then take the diagonal entries of this matrix to be the eigenvalues of A .

The basic idea is to note that we can diagonalize symmetric 2×2 matrices. Using this we choose the matrix S_{i+1} so that it makes some off-diagonal entry of $S_i^T \cdots S_1^T A S_1 \cdots S_i$ equal to 0. Unfortunately this will usually make some off-diagonal entries non-zero, when they were already zero. This will make us work harder, but will not prevent eventual success.

If M is a symmetric matrix then there is an orthogonal matrix L such that $L^T M L$ is diagonal; if M is 2×2 then we may assume that L has the form

$$\begin{pmatrix} c & -s \\ s & c \end{pmatrix}.$$

where $c^2 + s^2 = 1$. (We could, but do not, assume that $c \geq 0$.) Now suppose that A is a symmetric $n \times n$ matrix, that B is the leading principal 2×2 submatrix of A and that R is an orthogonal matrix such that $R^T B R$ is diagonal. Let S denote the matrix

$$\begin{pmatrix} R & 0 \\ 0 & I_{n-2} \end{pmatrix}.$$

Then

$$B = S^T A S$$

is similar to A and $B_{1,2} = 0$. In general, if $A_{i,j} \neq 0$ then there is an orthogonal matrix S such that the only non-zero off-diagonal entries of S are the ij and ji entries, and $(S^T A S)_{i,j} = 0$. We call S a *Givens rotation*.

How does this help us. If A and B are $n \times n$ matrices, define

$$\langle A, B \rangle = \text{tr } A B^T.$$

Then $\|A\|^2$ is the sum of the squares of the entries of A and, if L is orthogonal,

$$\|L^T A L\|^2 = \text{tr}(L^T A^T L L^T A^T L) = \text{tr}(L^T L A A^T L) = \text{tr}(A A^T) = \|A\|^2.$$

Let $\text{sqr}(A)$ denote the sum of the squares of the off-diagonal entries of A . Note that, in passing from A to $S^T A S$, the only diagonal entries that change are the ii - and jj -entries and that the sum of the squares of these two entries increases by $2(A_{i,j})^2$. It follows that, if S and A are as above,

$$\text{sqr}(S^T A S) = \text{sqr}(A) - 2(A_{i,j})^2.$$

If $\text{sqr}(A) = c$, there are indices i and j such that

$$(A_{i,j})^2 \geq \frac{c}{n(n-1)}$$

and hence there is a Givens rotation S such that

$$\text{sqa}(S^T AS) \leq c \left(1 - \frac{2}{n(n-1)} \right).$$

This implies that, by successively applying Givens rotations, we can form a matrix M , orthogonally similar to A and such that $\text{sqa}(M)$ is as small as we like. The diagonal entries of M will be the eigenvalues of A .

5.9 Jacobi: An Example

By way of example, suppose that

$$A = \begin{pmatrix} 1 & 0.5 & 0.3333 \\ 0.5 & 0.3333 & 0.25 \\ 0.3333 & 0.25 & 0.2 \end{pmatrix}.$$

Then Jacobi's method runs through the following iterations.

$$\begin{pmatrix} c \\ s \end{pmatrix} = \begin{pmatrix} -0.47185 \\ 0.88167 \end{pmatrix}; \quad [1,2] \rightarrow \begin{pmatrix} 0.065741 & 0 & 0.063132 \\ 0 & 1.2675 & -0.41185 \\ 0.063132 & -0.41185 & 0.2 \end{pmatrix}$$

$$\begin{pmatrix} c \\ s \end{pmatrix} = \begin{pmatrix} -0.32269 \\ -0.94650 \end{pmatrix}; \quad [2,3] \rightarrow \begin{pmatrix} 0.06574 & -0.05975 & -0.02037 \\ -0.05975 & 0.05958 & 0 \\ -0.02037 & 0 & 1.40801 \end{pmatrix}$$

$$\begin{pmatrix} c \\ s \end{pmatrix} = \begin{pmatrix} -0.68867 \\ -0.72507 \end{pmatrix}; \quad [1,2] \rightarrow \begin{pmatrix} 0.002829 & 0 & 0.01403 \\ 0 & 0.12250 & -0.01477 \\ 0.01403 & -0.01477 & 1.40801 \end{pmatrix}$$

$$\begin{pmatrix} c \\ s \end{pmatrix} = \begin{pmatrix} -0.99993 \\ -0.01149 \end{pmatrix}; \quad [2,3] \rightarrow \begin{pmatrix} 0.00283 & -0.00016 & -0.01403 \\ -0.00016 & 0.12232 & 0 \\ -0.01403 & 0 & 1.40818 \end{pmatrix}$$

$$\begin{pmatrix} c \\ s \end{pmatrix} = \begin{pmatrix} -0.99995 \\ -0.00998 \end{pmatrix}; \quad [1,3] \rightarrow \begin{pmatrix} 0.00269 & 0.00016 & 0 \\ 0.00016 & 0.12233 & 0 \\ 0 & 0 & 1.40832 \end{pmatrix}$$

$$\begin{pmatrix} c \\ s \end{pmatrix} = \begin{pmatrix} -1.0 \\ 0.00135 \end{pmatrix}; \quad [1,2] \rightarrow \begin{pmatrix} 0.00268 & 0 & 0 \\ 0 & 0.122327 & 0 \\ 0 & 0 & 1.40832 \end{pmatrix}$$

Here the diagonal entries are the eigenvalues of A , and further iterations do not change them.

- (1) Suppose that $\langle Av, Aw \rangle = 0$ whenever $\langle v, w \rangle = 0$. Prove, or disprove, that A is a scalar multiple of an orthogonal matrix.
- (2) Suppose $Q^2 = Q$ and $Q = Q^T$. Show that $I - 2Q$ is a symmetric orthogonal matrix, and explain the connection to reflections.
- (3) Prove that an involution is symmetric if and only if it is orthogonal.
- (4) Show that each involution has the form $I - 2P$, for some idempotent P .
- (5) Show that, if A and A^{-1} are similar, there is an involution T such that $TAT = A^{-1}$.

Chapter 6

Spectral Decomposition

6.1 Self-Adjoint Operators

The spectral decomposition of an operator is a more concrete form of diagonalizability. It is most useful when the operator is self-adjoint, so we confine ourselves to that case.

Suppose S is an operator on the inner product space V and that the minimal polynomial ψ of S is given by

$$\psi(t) = \prod_{i=1}^k (t - \theta_i),$$

where the zeros θ_i are distinct. (Thus A is diagonalizable.) By the primary decomposition theorem (Theorem 2.6.1), there are polynomials p_i such that

$$I = \sum_{i=1}^k p_i(S), \tag{6.1.1}$$

where

- (a) $p_i(S)$ is idempotent,
- (b) $p_i(S)p_j(S) = 0$ if $i \neq j$, and
- (c) S acts on $\text{col}(p_i(S))$ as multiplication by θ_i .

Assume $E_i := p_i(S)$. Then $SE_i = \theta_i E_i$ and, by (6.1.1),

$$S = \sum_i \theta_i E_i. \tag{6.1.2}$$

Equation (6.1.2) is known as the *spectral decomposition* of S .

One consequence of the spectral decomposition is that

$$S^n = \sum_i \theta_i^n E_i;$$

this can provide a simple way to compute powers of S .

If S is self-adjoint, then the operators E_i are self-adjoint, because each E_i is a polynomial in S .

We now offer a matrix view of the spectral decomposition. If A is a diagonalizable $n \times n$ matrix, then \mathbb{F}^n has a basis consisting of eigenvectors for A . Let L be the matrix with these eigenvectors as its columns. Then L is an invertible matrix and there is a diagonal matrix D such that

$$AL = LD.$$

It follows that

$$A = LDL^T.$$

We can write D as a sum of 01-diagonal matrices D_i :

$$D = \sum_i \theta_i D_i,$$

where $\theta_1, \dots, \theta_m$ are the distinct eigenvalues of A and $\sum_i D_i = I$. Accordingly

$$A = \sum_i \theta_i LD_i L^{-1}.$$

It is easy to verify that

$$(LD_i L^{-1})^2 = LD_i L^{-1}$$

and, if $i \neq j$, then $D_i D_j = 0$ and

$$LD_i L^{-1} LD_j L^{-1} = LD_i D_j L^{-1} = 0$$

6.2 Commutative Algebras

Two idempotents E and F are *orthogonal* if $EF = 0$. For example, if E is an idempotent, then E and $I - E$ are orthogonal idempotents. We can define a partial ordering on the idempotents of a commutative algebra \mathcal{A} as follows. If

E and F are idempotents in \mathcal{A} , we declare that $E \leq F$ if $FE = E$. This relation is reflexive, antisymmetric and transitive; therefore it is a partial order. A *minimal idempotent* is a minimal element of the set of non-zero idempotents, relative to this order. If E and F are idempotents, then $EF \leq E, F$. It follows that if E and F are minimal, then they are orthogonal.

6.2.1 Theorem. *Let \mathcal{B} be a commutative matrix algebra with identity over an algebraically closed field. Assume that if $N \in \mathcal{B}$ and $N^2 = 0$, then $N = 0$. Then \mathcal{B} has a basis of pairwise orthogonal idempotents.*

Proof. As a first step, we show that each element of \mathcal{B} is a linear combination of idempotents.

Assume the matrices in \mathcal{B} have order $n \times n$. Suppose $A \in \mathcal{B}$ and let $\psi(t) = \prod_{i=1}^k (t - \theta_i)^{m_i}$ be its minimal polynomial. There are idempotents E_i , summing to I , such that $\text{im}(E_i)$ is the root space associated with θ_i , and \mathbb{F}^n is the direct sum of these root spaces.

Further, the minimal polynomial of A on $\text{im}(E_i)$ is $(t - \theta_i)^{m_i}$, and hence we have

$$0 = (A - \theta_i I)^{m_i} E_i = ((A - \theta_i I) E_i)^{m_i}.$$

If $m_i > 1$, we set $k = \lfloor (m_i + 1)2 \rfloor$ and $N = ((A - \theta_i I) E_i)^k$. Then $N \neq 0$ but $N^2 = 0$. We conclude that zeros of the minimal polynomial of A are simple. We also see that $\text{im}(E_i)$ is an eigenspace for A and as $I = \sum_i E_i$ it follows that

$$A = AI = \sum_i A E_i = \sum_i \theta_i E_i.$$

Therefore A is a linear combination of idempotents belonging to \mathcal{B} , and it follows that \mathcal{B} is spanned by idempotents.

The problem that remains is to show that minimal idempotents exist. Suppose E and F are distinct idempotents and $E \leq F$. Then

$$F(I - E) = F - E \neq 0$$

but $E(I - E) = 0$. Hence the column space of E must be a proper subspace of the column space of F . Therefore if E_1, \dots, E_m are distinct idempotents and

$$E_1 \leq \dots \leq E_m$$

then $m \leq n + 1$. We conclude that minimal idempotents exist.

Now we prove that each idempotent is a sum of minimal idempotents. Suppose F is an idempotent and E is a minimal idempotent. If $EF \neq 0$, then $EF \leq E$ and therefore $EF = E$. This also shows that distinct minimal idempotents are orthogonal. Let F_0 be the sum of the distinct minimal idempotents E such that $E \leq F$. Then F_0 is an idempotent. If $F_0 \neq F$ then $F - F_0$ is an idempotent and so there is a minimal idempotent below it, which contradicts our choice of F_0 . We conclude that \mathcal{B} is spanned by minimal idempotents. \square

A matrix N is nilpotent if $N^k = 0$ for some k . Theorem 6.2.1 asserts that a commutative matrix algebra with identity has a basis of orthogonal idempotents if there are no non-zero nilpotent matrices in it. Since a non-zero linear combination of pairwise orthogonal idempotents cannot be nilpotent, this condition is necessary too. A commutative algebra is *semisimple* if it contains no non-zero nilpotent elements.

6.3 Normal Operators

An operator A on an inner product space is *normal* if $AA^* = A^*A$. We consider examples. Clearly any self-adjoint operator is normal. Unitary operators are a second important class. If $A = L^*DL$ where D is diagonal and L is unitary, then

$$AA^* = L^*DLL^*\bar{D}L = L^*D\bar{D}L = L^*\bar{D}DL = A^*A$$

and so any matrix that is unitarily similar to a diagonal matrix is normal.

Exercise: determine which complex 2×2 matrices are normal.

Exercise: If H is normal, show that we can write it as $H = A + iB$, where A and B are Hermitian and commute.

6.3.1 Theorem. Suppose \mathcal{A} is a commutative subalgebra of $\text{Mat}_{\nu \times \nu}(\mathbb{C})$ that is closed under conjugate transpose and contains the identity. Then \mathcal{A} has a basis of matrix idempotents E_0, \dots, E_d such that

(a) $E_i E_j = \delta_{i,j} E_i$.

(b) The columns of E_i are eigenvectors for each matrix in \mathcal{A} .

(c) $\sum_{i=0}^d E_i = I$.

(d) $E_i^* = E_i$. \square

Proof. Suppose $N \in \mathcal{A}$ and $N^2 = 0$. Then

$$0 = (N^*)^2 N^2 = (N^* N)^2$$

and hence

$$0 = \operatorname{tr}((N^* N)^2) = \operatorname{tr}((N^* N)^* (N^* N)).$$

If $H := N^* N$, then $\operatorname{tr}(H^* H) = 0$ if and only if $H = 0$, so we deduce that $N^* N = 0$. But then $\operatorname{tr}(N^* N) = 0$ and therefore $N = 0$. Hence \mathcal{A} satisfies the hypotheses of 6.2.1, and therefore it has a basis that consists of pairwise orthogonal idempotents.

We show that the idempotents E_i are Hermitian. Since \mathcal{A} is closed under transpose and complex conjugation, $E_i^* \in \mathcal{A}$. Therefore there are scalars a_0, \dots, a_d such that

$$E_i^* = \sum_j a_j E_j$$

and so

$$E_i^* E_i = f_i E_i.$$

Since $\operatorname{tr}(E_i^* E_i) > 0$ and $\operatorname{tr}(E_j) > 0$, it follows that $f_i \neq 0$. But E_i^* is a minimal idempotent, and therefore $f_j = 0$ if $j \neq i$. This implies that E_i^* is a scalar multiple of E_i , but $\operatorname{tr}(E_i) = \operatorname{tr}(E_i^*)$, and therefore $E_i^* = E_i$. \square

6.3.2 Theorem. *If A is normal, then A is unitarily similar to a diagonal matrix.*

Chapter 7

Norms

7.1 Convexity

We work over \mathbb{R} . We say that a vector v is an *affine combination* of vectors x_1, \dots, x_n if

$$v = \sum_i a_i x_i$$

and $\sum a_i = 1$. An affine combination is *proper* if it has at least two non-zero coefficients. The set of all affine combinations of a set of vectors is the *affine hull* of the set. The affine hull of x is x itself. The affine hull of $\{x, y\}$ (where $x \neq y$) is

$$\{tx + (1-t)y : t \in \mathbb{R}\}.$$

Geometrically this set is the unique line passing through the points represented by x and y . Note that this line contains 0 if and only if x and y are linearly dependent.

If U is a subspace of V and then a *coset* of U is a set of the form

$$\{a + u : u \in U\},$$

for some a in V .

7.1.1 Lemma. *The affine hull of a set of vectors $\{x_1, \dots, x_m\}$ is a coset of the subspace spanned by $x_2 - x_1, \dots, x_m - x_1$.* \square

An *affine subspace* is a set S that is closed under affine combinations.

We say that vectors x_1, \dots, x_m are *affinely dependent* if there are scalars a_i , not all zero, such that

$$\sum_i a_i = 0, \quad \sum_i a_i x_i = 0.$$

If a set is not affinely dependent, it is *affinely independent*. Note that any single vector, including the zero vector, is affinely independent.

A vector v is a *convex combination* of vectors x_1, \dots, x_m if there are scalars a_1, \dots, a_m such that

$$\sum_i a_i = 1, \quad a_i \geq 0 \quad (i = 1, \dots, m)$$

and

$$v = \sum a_i x_i.$$

Thus a convex combination is a non-negative affine combination. A convex combination is *proper* if it has at least two non-zero coefficients. The *convex hull* of a subset S is the set of all convex combinations of elements of S . A set S is *convex* if any convex combinations of its elements is contained in S , that is, if S is equal to its convex hull.

The convex hull of two distinct vectors consists of the line segment that joins them. Hence a set S is convex if, whenever x and y belong to S , so do all points on the line segment joining them. We also see that the intersection of two convex sets is convex.

A real-valued function f on \mathbb{R}^n is *convex* if

$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y), \quad 0 \leq t \leq 1.$$

- (1) If $a \in \mathbb{R}^n$, show that $f(x) := \exp(a^T x)$ is a convex function.
- (2) Show that set of positive semidefinite matrices is the convex hull of the matrices with rank 1.
- (3) Suppose $a_i \geq 0$ and $\sum_i a_i = 1$. If f is convex, prove that

$$f\left(\sum_i a_i x_i\right) \leq \sum_i a_i f(x_i).$$

- (4) Use the result of the previous exercise with $f(x) = x^p$ ($p > 1$) to show that

$$\sum_i |x_i y_i| \leq (|x_i|^p)^{1/p} (|y_i|^q)^{1/q},$$

where $1/p + 1/q = 1$. (This is Hölder's inequality.)

7.2 Extreme Points

Let C be a convex set. A point x in C is extreme if it cannot be expressed as the convex combination of points in $C \setminus x$. The extreme points of a line segment are its endpoints. Suppose C is convex and $x \in C$. Let ℓ be a line through x . Then $\ell \cap C$ is a line segment. The interior points of this line segment are not extreme. A closed convex set is the convex hull of its extreme points. We will not prove this, but we consider two cases that will be useful.

7.2.1 Lemma. *Let S be the set of vectors x in \mathbb{R}^n such that $|x_i| \leq 1$ for all i . Then S is the convex hull of the vectors with all entries ± 1 .*

Proof. It is easy to verify that S is convex, we leave this as an exercise. We show that it is the convex hull of the ± 1 -vectors.

We prove this by induction on n , asserting that it is trivial when $n = 1$. Assume $v \in S$ and that $v_1 = 1$. Let v' be the vector we get by deleting the first entry of v . Then v' lies in the set of vectors x in \mathbb{R}^{n-1} such that $|x_i| \leq 1$, and so by induction it is a convex combination of the ± 1 -vectors in \mathbb{R}^{n-1} . It follows that v is a convex combination of those ± 1 -vectors in \mathbb{R}^n with first entry equal to 1. If $v_1 = -1$, then $-v'$ is a convex combination of ± 1 -vectors x_1, \dots, x_m , and so v' is a convex combination of the vectors $-x_1, \dots, -x_m$, but these are ± 1 -vectors too. It follows that if $|v_i| = 1$, then v is a convex combination of ± 1 -vectors.

Now suppose that $|v_i| < 1$ for all i . Let v^+ be the vector such that $(v^+)_i = 1$ if $v_i \geq 0$ and $(v^+)_i = -1$ if $v_i < 0$. Then

$$((1-t)v^+ + tv)_i = \begin{cases} 1-t+tv_i, & \text{if } v_i \geq 0; \\ t-1+tv_i, & \text{otherwise.} \end{cases}$$

from which we eventually deduce that $w = (1-t)v^+ + tv \in S$ provided

$$0 \leq t \leq \frac{2}{1-|v_i|}.$$

Choose t so that $t = 2/(1-|v_j|)$ for some j . Then $|w_j| = 1$, and therefore v is a convex combination of v^+ and w . Since $|w_j| = 1$, it is the convex combination of ± 1 -vectors, and therefore v is too. \square

7.2.2 Lemma. *Let S be the set of vectors x such that*

$$\sum_i |x_i| \leq 1.$$

Then S is the convex hull of the vectors $\pm e_i$ for $i = 1, \dots, n$. \square

- (1) Show that if x is a proper convex combination of points from C , it is the proper convex combination of two points.
- (2) Let C be a convex set and let f be a convex function. If the point x_0 in C maximizes the value of f , show that it is an extreme point.
- (3) Prove (??).

7.3 Norms

Let V be a vector space over \mathbb{F} , where \mathbb{F} is \mathbb{R} or \mathbb{C} . A *norm* on V is a function from V to \mathbb{R} , whose value on x is written $\|x\|$, such that

- (1) $\|x\| \geq 0$ and $\|x\| = 0$ if and only if $x = 0$.
- (2) If $c \in \mathbb{R}$, then $\|cx\| = |c| \|x\|$.
- (3) If $\|x + y\| \leq \|x\| + \|y\|$.

The third axiom is called the *triangle inequality*. It implies that any norm is a convex function on V . The set

$$\{x : \|x\| \leq 1\}$$

is called the *unit ball* of the norm, but it need not be very round.

We consider some examples over \mathbb{R} . If we have an inner product on V , then we can define a norm by

$$\|x\| := \sqrt{\langle x, x \rangle}$$

The only difficulty here is to verify the triangle inequality. We note that

$$\|x + ty\|^2 = \langle x + ty, x + ty \rangle = \langle x, x \rangle + 2\langle x, y \rangle t + \langle y, y \rangle t^2.$$

This is a quadratic in t which is non-negative for all t , and consequently

$$\langle x, y \rangle^2 - \langle x, x \rangle \langle y, y \rangle \leq 0,$$

which is usually called the Cauchy-Schwarz inequality. It follows that

$$\begin{aligned} \langle x, x \rangle + 2\langle x, y \rangle t + \langle y, y \rangle t^2 &\leq \langle x, x \rangle + 2\langle x, y \rangle t + \langle y, y \rangle t^2 \\ &\leq \langle x, x \rangle + 2\|x\| \|y\| t + \langle y, y \rangle t^2 \\ &= (\|x\| + t\|y\|)^2. \end{aligned}$$

We conclude that $\|x\| + \|ty\| \leq \|x + ty\|$, which yields the triangle inequality.

If our inner product is the dot product our norm is the usual Euclidean norm or ℓ_2 -norm and is denoted by $\|\cdot\|_2$ or, sometimes, by $\|\cdot\|$. The unit ball for the Euclidean norm is the unit ball.

If $\langle \cdot, \cdot \rangle$ is a complex inner product, the function

$$\sqrt{\langle x, x \rangle}$$

is a norm. Note that $\langle x, x \rangle$ is guaranteed to be real and non-negative.

Once we have a norm, we can declare that a sequence x_0, x_1, \dots of vectors converges to x if the sequence of real numbers

$$\|x - x_0\|, \|x - x_1\|, \dots$$

converges to 0. It is a somewhat surprising fact that if a sequence of vectors in a finite-dimensional vector space converges with respect to one norm, then it converges with respect to all. (This is false if the dimension is infinite, as the exercises show.)

- (1) Prove that a norm is a convex function.
- (2) Let $V = C[0, 1]$, the space of continuous functions on the interval $[0, 1]$. If $f \in V$, let $\|f\|$ be the norm associated with the inner product

$$\langle f, g \rangle := \int_0^1 f(x)g(x) dx$$

and let $\|f\|_\infty$ be the norm defined by

$$\|f\|_\infty = \max\{f(x) : x \in [0, 1]\}.$$

(You may prove that this is a norm.) Define

$$g_r(x) := (4x(1-x))^r.$$

Prove that $\|g_r\| \rightarrow 0$ as $r \rightarrow \infty$, but $\|g_r\|_\infty = 1$ for all r .

7.4 Dual Norms

We introduce two further norms. We define $\|x\|_1$ by

$$\|x\|_1 := \max_i \sum_i |x_i|$$

and $\|x\|_\infty$ by

$$\|x\|_\infty := \max_i |x_i|.$$

These are known respectively as the ℓ_1 and ℓ_∞ -norms on \mathbb{R}^n . As we saw in the previous section, the unit ball for the ℓ_1 -norm is the convex hull of the vectors $\pm e_i$ and the unit ball for the ℓ_∞ -norm is the convex hull of the ± 1 vectors. (These definitions work over both \mathbb{R} and \mathbb{C} , we will only use them over \mathbb{R} though.)

If $\|\cdot\|$ is a norm, we define the *dual norm* $\|\cdot\|^*$ by

$$\|a\|^* := \max_{\|x\|=1} x^T a.$$

We leave the proof that this is a norm as an exercise. As another exercise, we leave you to prove that $\|x\|^{**} = \|x\|$, for any x .

By way of example, we determine the dual of the ℓ_∞ -norm. Our problem is compute the maximum value of the function $x^T a$ over the vectors x in the unit ball of the ℓ_∞ -norm. This is linear in x , and hence convex; therefore its maximum value occurs at an extreme point of this ball. By Lemma 7.2.1, the extreme points are the ± 1 -vectors and hence $\|a\|_\infty^*$ is equal to the maximum value of $x^T a$, as x ranges over the set of ± 1 -vectors. Clearly this maximum is realized when $x_i a_i > 0$ for each i , and therefore

$$\|a\|_\infty^* = \sum_i |a_i| = \|a\|_1.$$

- (1) Let V be the Euclidean space \mathbb{R}^n . Determine the largest C and the smallest D such that

$$C\|x\|_\infty \leq \|x\| \leq D\|x\|_\infty.$$

- (2) If the function $\|\cdot\|^*$ is defined on \mathbb{R}^n by

$$\|y\|^* = \max_{\|x\|=1} x^T y,$$

show that it is a norm.

- (3) Prove that $\|x\|^{**} = \|x\|$, for any x .
- (4) Prove that $y^T x \leq \|x\| \|y\|^*$, and show that this bound is tight.
- (5) Show that the ℓ_1 -norm is dual to the ℓ_∞ -norm, and vice versa.

7.5 Matrix Norms

Let \mathcal{B} be an algebra over the reals. A norm on \mathcal{B} is a function $\|\cdot\|$ from \mathcal{B} to \mathbb{R} that is norm, when we view \mathcal{B} as a vector space, and in addition satisfies:

$$\|AB\| \leq \|A\| \|B\|.$$

If $\|\cdot\|$ is a norm on a vector space V then the unit ball

$$\{x \in V : \|x\| \leq 1\}$$

is a closed convex set. If $\|\cdot\|$ is a norm on an algebra then the unit ball must be closed under multiplication, hence forms a semigroup.

Now suppose $\|\cdot\|$ is a norm on $L(V)$, viewed as a vector space. The unit ball is compact and so, if $A \in L(V)$ then there is a constant γ_A such that, if $\|X\| \leq 1$,

$$\|AX\| \leq \gamma_A.$$

If we define γ to be the maximum value of γ_A , where $\|A\| \leq 1$, then

$$\|AB\| = \|A\| \|B\| \gamma.$$

From this it follows that $\gamma^{-1}\|\cdot\|$ is a norm on $L(V)$, viewed as an algebra. We will refer to a norm on an algebra as an *operator norm* or *matrix norm*, according as the elements of our algebra are linear mappings or matrices.

Let V be a normed vector space, with norm $\|\cdot\|$. If T is an endomorphism of V , we define the *induced norm* of T by

$$\|T\| = \max\{\|Tx\| : \|x\| = 1\}.$$

Equivalently, it is the maximum value of $\|Tx\|/\|x\|$, for all non-zero vectors x in V . It is straightforward to verify that this is a norm on $L(V)$, with the useful properties:

$$\|Tx\| \leq \|T\| \|x\|$$

and

$$\|ST\| \leq \|S\| \|T\|.$$

Unless explicitly stated otherwise, we use the same symbol to denote a norm on \mathbb{R}^n and the norm it induces on $n \times n$ matrices. If $\|\cdot\|$ is an induced norm, then $\|I\| = 1$.

If $\|\cdot\|$ is an induced norm then for any matrix A and vector x , we have the very useful inequality:

$$\|Ax\| \leq \|A\| \|x\|.$$

If $\|\cdot\|_a$ and $\|\cdot\|_b$ are any two norms on a vector space, we say that $\|\cdot\|_b$ *dominates* $\|\cdot\|_a$ if, for all v in V ,

$$\|v\|_a \leq \|v\|_b.$$

A norm is *minimal* if it does not dominate any other norm. Generally minimal norms are more useful than general norms.

7.5.1 Lemma. *Every matrix norm dominates an induced norm.*

Proof. Suppose $\|\cdot\|$ is a matrix norm. We use this to construct a norm on \mathbb{R}^n whose induced norm is dominated by $\|\cdot\|$.

Let a be a fixed non-zero vector in \mathbb{R}^n . We define $\|\cdot\|_a$ by

$$\|b\|_a := \|ba^T\|.$$

Then

$$\|Ax\|_a = \|Axa^T\| \leq \|A\| \|x\|_a$$

and the matrix norm induced by $\|\cdot\|_a$ is dominated by $\|\cdot\|$. \square

7.5.2 Theorem. *Let $\|\cdot\|$ be a norm on \mathbb{R}^n with dual norm $\|\cdot\|^*$. If A is a square matrix then $\|A\|^* = \|A^T\|$.*

Proof. We have

$$\|Ax\|^* = \max_{\|y\|=1} y^T Ax = \max_{\|y\|=1} x^T A^T y$$

and so

$$\|A\|^* = \max_{\|x\|^*=1} \max_{\|y\|=1} x^T A^T y.$$

Now

$$\max_{\|x\|^*=1} x^T b = \|b\|^{**} = \|b\|$$

and consequently

$$\|A\|^* = \max_{\|y\|=1} \|A^T y\| = \|A^T\|.$$

In the sequel any norm we use on matrices will be a matrix norm. If $(A_n)_{n \geq 0}$ is a sequence of matrices and we write that $A_n \rightarrow 0$, we mean that $\|A_n\| \rightarrow 0$, for some norm $\|\cdot\|$.

- (1) Let $\|\cdot\|$ be a norm on \mathbb{R}^n , and let $\|\cdot\|$ also denote the induced matrix norm. Prove that $\|ab^T\| = \|a\| \|b\|^*$ and hence that $b^T a \leq \|ab^T\|$.
- (2) Prove that if $n \geq 1$, then $\|A^n\|^{1/n} \leq \|A\|$.

7.6 Examples

The *Euclidean* or *trace* norm of a matrix is the norm associated with the inner product

$$\langle A, B \rangle := \text{tr } A^T B.$$

We denote this norm by $\|\cdot\|_2$ or, sometimes, by $\|\cdot\|$. Note that $\|A\|_2^2$ is the sum of the squares of the entries of A . We have

$$\begin{aligned} \|AB\|_2^2 &= \sum_{i,j} \left| \sum_r A_{i,r} B_{r,j} \right|^2 \leq \sum_{i,j} \left(\sum_r |A_{i,r}|^2 \right) \left(\sum_r |B_{r,j}|^2 \right) \\ &= \left(\sum_{i,r} |A_{i,r}|^2 \right) \left(\sum_{r,j} |B_{r,j}|^2 \right) \\ &= \|A\|_2^2 \|B\|_2^2. \end{aligned}$$

We have $\|I_n\| = n$ and so the trace norm is not an induced norm.

We turn next to induced matrix norms. First we note that

$$\|Ax\|_2^2 = (Ax)^T Ax = x^T A^T Ax$$

and therefore

$$\max_{\|x\|_2=1} \|Ax\|_2$$

is equal to $\sqrt{\rho}$, where ρ is the largest eigenvalue of $A^T A$. (But since we have not discussed eigenvalues at any length yet, we defer any further discussion.)

Both of above norms have the useful property that, if Q is orthogonal, then $\|QA\| = \|A\|$.

7.6.1 Lemma. *Let A be a square matrix. Then*

$$\|A\|_\infty = \max_i \|e_i^T A\|_1.$$

Proof. The function $x \mapsto \|Ax\|_\infty$ is convex and hence realizes its maximum at an extreme point of the unit ball relative to the ℓ_∞ norm. These extreme points are the ± 1 -vectors. If x is a ± 1 -vector then

$$|(Ax)_i| = \left| \sum_j A_{i,j} x_j \right| \leq \sum_j |A_{i,j} x_j| \leq \sum_j |A_{i,j}| = \|e_i^T A\|_1.$$

Further, equality holds throughout if we choose x so that $A_{i,j} x_j \geq 0$. This proves the lemma. \square

7.6.2 Lemma. *Let A be a square matrix. Then*

$$\|A\|_1 = \max_i \|Ae_i\|_1.$$

Proof. Since ℓ_1 and ℓ_∞ -norms are dual, we can apply Theorem 7.5.2 to the previous lemma, concluding that

$$\|A\|_1 = \|A^T\|_\infty = \max_i \|e_i^T A^T\|_1 = \max_i \|Ae_i\|_1. \quad \square$$

(1) If $\|\cdot\|$ is the trace norm or the induced ℓ_2 -norm, and Q is an orthogonal matrix, show that $\|QA\| = \|A\|$

7.7 Matrix Functions

We say a matrix is a function of a variable t if each element of the matrix is. This makes sense over any field, but here we work over \mathbb{R} or \mathbb{C} . If the matrix $A(t)$ is a function of t then

$$\frac{d}{dt} A(t)$$

is the matrix we get by differentiating each entry of $A(t)$ with respect to t .

As an example, we consider the differential equation

$$f'' + af' + b = 0. \quad (7.7.1)$$

This is equivalent to the following pair of equations:

$$\begin{aligned}\frac{d}{dt}f' &= -af' - b, \\ \frac{d}{dt}f &= f',\end{aligned}$$

which we can rewrite as

$$\frac{d}{dt}\begin{pmatrix} f' \\ f \end{pmatrix} = \begin{pmatrix} -a & -b \\ 1 & 0 \end{pmatrix} \begin{pmatrix} f' \\ f \end{pmatrix}.$$

We can solve this using the matrix exponential.

For any square matrix A we define

$$\exp(tA) := \sum_{n=0}^{\infty} \frac{t^n}{n!} A^n.$$

But we need to see that this makes sense. We have

$$\|A^n\|_{\infty} \leq \|A\|_{\infty}^n$$

and so, if $a := \|A\|_{\infty}$, each entry of A^n is bounded in absolute value by a^n . Therefore each entry of

$$\sum_{n=0}^m \frac{t^n}{n!} A^n$$

converges as $m \rightarrow \infty$, for any value of t . Moreover we are entitled to differentiate the series term-by-term, with the result that

$$\frac{d}{dt} \exp(tA) = \sum_{n=1}^{\infty} \frac{t^{n-1}}{(n-1)!} A^n = A \exp(tA).$$

Now define the vector $F(t)$ by

$$F(t) = \begin{pmatrix} f' \\ f \end{pmatrix}$$

and suppose

$$A := \begin{pmatrix} -a & -b \\ 1 & 0 \end{pmatrix}.$$

Then (7.7.1) becomes

$$\frac{d}{dt}F(t) = AF(t)$$

and it is easy to see that this has the solution

$$F(t) = \exp(tA)F(0).$$

Although this method of solving differential equations is very important, it is of limited use as a tool for solving particular equations. It is computationally difficult to compute $\exp(A)$ because, even though

$$\frac{1}{n!}A^n \rightarrow 0$$

as $n \rightarrow \infty$, for moderate values of n this ratio can be very large. The difficulty is essentially the same as attempting to compute $\exp(100)$ using the power series for the exponential.

(1) Show that

$$\exp t(A + B) = \exp(tA) \exp(tB)$$

if and only if $AB - BA = 0$.

(2) If S is skew symmetric, show that $\exp(S)$ is orthogonal.

(3) If

$$H := \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix},$$

show that $\exp(\pi H) = -I$.

7.8 Powers

We have seen the exponential series in a matrix A is well-defined and useful. We will find useful to consider other power series with matrix arguments. Our next result provides a basic tool.

7.8.1 Lemma. *If A is a non-zero matrix and $\|\cdot\|$ is a matrix norm, then the sequence $\|A^n\|^{1/n}$ converges to a limit ρ . Further $\rho \leq \|A^n\|^{1/n}$ for all n .*

Proof. By way of abbreviation, let $f(n) = \|A^n\|^{1/n}$. Note first that

$$\|A^{km}\| \leq \|A^m\|^k,$$

and therefore $f(km) \leq f(m)$. Assume $n = km + \ell$, where $0 \leq \ell < m$. Then

$$f(km + \ell) \leq f(km)^{\frac{km}{km+\ell}} f(\ell)^{\frac{\ell}{km+\ell}} \leq f(m)^{\frac{km}{km+\ell}} f(\ell)^{\frac{\ell}{km+\ell}}$$

Given $\epsilon > 0$ and fixed m , it follows that for all but finitely many n , we have

$$f(n) \leq (1 + \epsilon)f(m).$$

We say that $f(m)$ is a *record* for f if, when $k < m$,

$$f(m) < f(k).$$

Consider the sequence of records for f . If it is finite, let ρ denote its last member. If it is not finite, then it is a strictly decreasing sequence, bounded below by 0 and therefore it has a limit, which we denote by ρ . From the previous paragraph it follows that if $\epsilon > 0$, then $f(n) \leq (1 + \epsilon)\rho$ for all but finitely many values of n . Consequently the sequence $\|A^n\|^{1/n}$ converges to ρ and $\rho \leq \|A^n\|^{1/n}$ for all n . \square

This lemma does not guarantee that $\rho^{-n}A^n$ converges. For example, if

$$A = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$$

then

$$A^n = \begin{pmatrix} \cos n\theta & -\sin n\theta \\ \sin n\theta & \cos n\theta \end{pmatrix}$$

and, using the trace norm

$$\|A^n\| = \|A\| = 2.$$

Therefore $\|A^n\|^{1/n} = 1$ but, nonetheless, the sequence $(A^n)_{n \geq 0}$ does not converge except in special cases.

The quantity

$$\lim_{n \rightarrow \infty} \|A^n\|^{1/n}$$

is known as the *spectral radius* of A .

We want to work with the geometric series

$$\sum_{r \geq 0} t^r A^r.$$

7.8.2 Lemma. *The series $\sum_{r \geq 0} t^r A^r$ converges if and only if $t^n A^n \rightarrow 0$ as $n \rightarrow \infty$. If it does converge, its limit is $(I - tA)^{-1}$.*

Proof. We have

$$(I - tA)(I + tA + \cdots + t^{n-1} A^{n-1}) = I - t^n A^n.$$

Suppose $I - tA$ is not invertible. Then there is a non-zero vector u such that $(I - tA)u = 0$. Therefore $tAu = u$ and $t^r A^r u = u$ for all r . So $t^n A^n$ does not converge to 0 and, since

$$(I + tA + \cdots + t^{n-1} A^{n-1})u = nu,$$

the series $\sum_{r \geq 0} t^r A^r$ does not converge.

Hence we may suppose that $I - tA$ is invertible and consequently

$$I + tA + \cdots + t^{n-1} A^{n-1} = (I - tA)^{-1}(I - t^n A^n).$$

The lemma follows immediately. \square

7.8.3 Corollary. *Let ρ be the spectral radius of A . The series $\sum_r t^r A^r$ converges (to $(I - tA)^{-1}$) if $|t| < \rho^{-1}$ and diverges if $|t| > \rho^{-1}$.*

Proof. We observe that $t^n A^n$ converges to 0 if and only if $\|t^n A^n\|$ does. By Lemma 7.8.1 we see that $t^n A^n \rightarrow 0$ if $|t| < \rho^{-1}$ and that it does not converge if $|t| > \rho$. \square

This result shows that ρ^{-1} is the radius of convergence of the series $\sum_n t^n A^n$.

7.9 Contractions

We call a linear map T a *contraction* relative to the norm $\|\cdot\|$ if $\|T^n\| \rightarrow 0$ as n increases. Our first result shows that being a contraction is independent of the norm we use.

7.9.1 Lemma. *The linear map T is a contraction if and only if its spectral radius is less than 1.*

Proof. Let ρ be the spectral radius of T . Let $\|\cdot\|$ be an operator norm, and suppose $\epsilon > 0$. By Lemma 7.8.1, for all sufficiently large values of n ,

$$\rho^n \leq \|T^n\| \leq (\rho + \epsilon)^n.$$

The result follows at once. \square

While this result has its uses, it does not provide an effective means of deciding if a particular map is a contraction. But contractions are important, and so we need effective ways of recognizing them. If there is an operator norm such that $\|T\| < 1$, then since

$$\|T^n\| \leq \|T\|^n,$$

it follows that T is a contraction. Our work in this section shows that, if T is a contraction, there is a norm $\|\cdot\|$ such that $\|T\| < 1$.

If B is a positive definite matrix then the bilinear form

$$\langle u, v \rangle = u^T B v$$

is an inner product, and $\sqrt{u^T B u}$ is a norm. (See Lemma 9.1.2.)

7.9.2 Lemma. *A matrix A is a contraction if and only if there is a positive definite matrix B such that $B - A^T B A$ is positive definite.*

Proof. Suppose first that B is positive definite and $B - A^T B A$ is positive definite. Then for any non-zero vector v ,

$$0 < v^T (B - A^T B A) v = v^T B v - v^T A^T B A v.$$

If $\|\cdot\|_B$ denotes the norm determined by B , this shows that, for any non-zero vector v ,

$$\|A v\|_B < \|v\|_B$$

and therefore $\|A\|_B < 1$.

To complete the proof, we show that if C is positive definite and the equation

$$X - A^T X A = C \tag{7.9.1}$$

has a positive definite solution X , then A is a contraction. If X satisfies (7.9.1), then

$$\begin{aligned} X &= C + A^T X A \\ A^T X A &= A^T C A + (A^T)^2 X A^2 \\ (A^T)^2 X A^2 &= (A^T)^2 C A^2 + (A^T)^3 X A^3, \end{aligned}$$

which leads us to conjecture (and you to prove, by summing enough of these equations) that

$$X - (C + A^T C A + \cdots + (A^T)^{n-1} C A^{n-1}) = (A^T)^n C A^n.$$

Since the right side of this identity goes to 0 as n increases, we conclude that

$$X = \sum_{r \geq 0} (A^T)^r C A^r$$

is a solution to (7.9.1). Because C is positive definite, $v^T C v > 0$ for all non-zero vectors v , and therefore

$$v^T (A^T)^r C A^r v > 0$$

for all non-zero vectors v . Consequently X is positive definite. □

Equation (7.9.1) is known as *Stein's equation*. It is a system of linear equations in the entries of X , and so can readily be solved. Since all we need of C is that it be positive definite, we may choose $C = I$. The proof of the lemma shows that if A is a contraction, then Stein's equation has a unique solution. Therefore we could determine if A is a contraction by solving $X - A^T X A = I$, and then testing whether the solution X is positive definite. (This can be decided by Cholesky factorization.)

- (1) If C is symmetric and $X - A^T X A = C$ has a solution, show that it has symmetric solution.
- (2) Read up on Kronecker products (in Corollary ??), and then show that, if A does not have distinct eigenvalues whose product is equal to 1, then $X - A^T X A = I$ has a solution.

7.10 Projections

We study subspaces and projections in \mathbb{R}^n ; our results extend to any inner product space. Suppose U is a k -dimensional subspace of \mathbb{R}^n , and let Y be an $n \times k$ matrix whose columns form a basis for U . The Gram-Schmidt algorithm implies that there is a $k \times k$ upper-triangular matrix P such that the columns of YP are orthogonal. As Y and YP have the same column space, it follows that the columns of YP form an orthonormal basis for U .

For our purposes we may as well assume that we chose Y so that $Y^T Y = I_k$, without further ado. If we define

$$P = Y^T Y$$

then we see that P is symmetric and

$$P^2 = Y^T Y Y^T Y = Y^T Y = P.$$

Hence P represents orthogonal projection onto its column space. As $\text{rk } P = \text{rk } Y = k$ and as the column space of P is contained in the column space of Y , it follows that the column space of P equals U . So P represents orthogonal projection onto U . One consequence of this is that the properties of the collection of k -dimensional subspaces of \mathbb{R}^n are mirrored by the properties of the $n \times n$ orthogonal projections with rank k .

Our projections are symmetric and there is a natural inner product on the space of symmetric matrices:

$$\langle A, B \rangle = \text{tr}(AB).$$

If $P_i = Y_i^T Y_i$ where Y_i is $n \times k$ and $Y_i^T Y_i = I_k$ then

$$\begin{aligned} \langle P_1, P_2 \rangle &= \text{tr}(Y_1 Y_1^T Y_2 Y_2^T) = \text{tr}(Y_2^T Y_1 Y_1^T Y_2) \\ &= \text{tr}((Y_1^T Y_2)^T (Y_1^T Y_2)) \\ &\geq 0. \end{aligned}$$

Further

$$\begin{aligned} \langle P_1 - P_2, P_1 - P_2 \rangle &= \text{tr}(P_1^2 - P_1 P_2 - P_2 P_1 + P_2^2) \\ &= \text{tr}(P_1 + P_2 - 2P_1 P_2) \\ &= 2k - 2\langle P_1, P_2 \rangle. \end{aligned}$$

Thus the value of $k - \text{tr}(P_1 P_2)$ can be viewed as a measure of how close the subspaces represented by P_1 and P_2 are.

If P and Q are projections defining two subspaces U and V of \mathbb{R}^n and x is a unit vector in \mathbb{R}^n then $\|Px - Qx\|$ is a measure of distance of U from V . Now

$$\|Px - Qx\|^2 = x^T (P - Q)^2 x,$$

whence all information of this sort is contained in the matrix $(P - Q)^2$. The maximum value over all unit vectors x of

$$\|Px - Qx\|^2 = x^T (P - Q)^2 x$$

is the largest eigenvalue of the (real symmetric) matrix $(P - Q)^2$. Our next result bounds this.

7.10.1 Lemma. *Let P and Q be projections. Then $\|Px - Qx\| \leq \|x\|$ and, if equality holds, $x = Px + Qx$ and $\langle Px, Qx \rangle = 0$.*

Proof. The vectors Px and $(I - P)x$ are orthogonal, so the points represented by the vectors 0 , Px and x are the vertices of a right-angled triangle with hypotenuse joining 0 to x . Thus (why??) they lie on the circle with this hypotenuse as a diameter. Similarly the vectors 0 , Qx and x form a second right-angled triangle, and also lie on a circle. Now, if two triangles in \mathbb{R}^n share a side then the distance between their third vertices is maximal when they lie in the same plane (and on opposite sides of their shared side). Hence $\|Px - Qx\| \leq \|x\|$; if equality holds then the two triangles are coplanar, the two circles coincide and Px and Qx must be diametrically opposed on the circle. Since the origin is on a circle with the line segment from Px to Qx as a diameter, Px , 0 and Qx form a right triangle and Px must be orthogonal to Qx . Further, 0 , Px , x and Qx form the vertices of a rectangle; by the parallelogram rule for addition of vectors in the plane, $x = Px + Qx$. \square

(1) Show that if P and Q are projections and $\text{rk } P = \text{rk } Q$, then $\text{tr}(P - Q)^3 = 0$.

(2) Show that $(P - Q)^2$ commutes with P and Q .

7.11 Contractions

In this section, we derive the characterization of contractions in terms of eigenvalues. If M is a square matrix, we use $\|M\|_1$ to denote the induced ℓ_1 norm of M —this equals the maximum value of the ℓ_1 -norms of the columns of M , as we saw in ??.

7.11.1 Theorem. *Let A be a square matrix. If $|\theta| < 1$ for all eigenvalues θ of A , then A is a contraction.*

Proof. As a first step, we prove the theorem when A is lower triangular. Suppose A is $n \times n$ and let D_t be the $n \times n$ diagonal matrix with $(D_t)_{i,i} = t^{i-1}$. Let Δ denote the diagonal matrix with $\Delta_{i,i} = A_{i,i}$. The ij -entry of $D_t^{-1}AD_t$ is $t^{j-i}A_{i,j}$ and so

$$\lim_{t \rightarrow \infty} D_t^{-1}AD_t = \Delta.$$

In particular, given $\epsilon > 0$, we can choose t large enough that $\|D_t^{-1}AD_t\|_1$ lies within ϵ of $\|\Delta\|_1$. Consequently, if $|\theta| < 1$ for each eigenvalue θ , then we can choose t so that $\|D_t^{-1}AD_t\|_1 < 1$.

This implies that

$$\|D_t^{-1}A^nD_t\|_1 \rightarrow 0$$

as $n \rightarrow \infty$. Since

$$\|A^n\|_1 = \|D_t^{-1}D_t^{-1}A^nD_tD_t\|_1 \leq \|D_t^{-1}\|_1\|D_t^{-1}A^nD_t\|_1\|D_t\|_1,$$

it follows that $\|A^n\|_1 \rightarrow 0$ as $n \rightarrow \infty$.

If A is not triangular, then $A = LTL^{-1}$, where T is triangular. Since

$$\|A^n\|_1 \leq \|L\|_1\|T^n\|_1\|L^{-1}\|_1$$

and

$$\|T^n\|_1 = \|L^{-1}A^nL\|_1 \leq \|L^{-1}\|_1\|A^n\|_1\|L\|_1,$$

we see that A is a contraction if and only if T is. To complete the proof, we recall that A and T have the same eigenvalues. \square

There is another proof of this result using root vectors.

7.11.2 Lemma. *Let A be an $n \times n$ matrix over \mathbb{C} , let θ be an eigenvalue of A and let v be a root vector for θ . If $|\theta| < 1$, then $A^m v \rightarrow 0$ as $m \rightarrow \infty$.*

Proof. Since v is a root vector for θ , we have $(A - \theta I)^n v = 0$. Then

$$A^m = (A - \theta I + \theta I)^m$$

and so using the binomial theorem, we find that

$$A^m v = \theta^{m-n+1} \left[\binom{m}{n-1} (A - \theta I)^{n-1} + \binom{m}{n-2} (A - \theta I)^{n-2} \theta + \cdots + \theta^{n-1} I \right] v.$$

Hence we have

$$A^m v = \theta^{m-n+1} P(m) v,$$

where $P(m)$ is a matrix whose entries are polynomials in m with degree at most $n - 1$. Since $|\theta| < 1$, it follows that

$$\theta^{m-n+1}P(m) \rightarrow 0$$

as $m \rightarrow \infty$. □

Now suppose A is an $n \times n$ matrix with all eigenvalues inside the unit circle. Since each vector in \mathbb{C}^n is a linear combination of root vectors, it follows that for any vector v ,

$$A^m v \rightarrow 0$$

as $m \rightarrow \infty$.

We have two methods now for determining if a square matrix A is a contraction. We can solve Stein's equation, as discussed in Section 7.9, or we can compute the spectral radius from the eigenvalues of A . This second alternative is useful if A is symmetric, or if A is real and its entries are positive.

7.12 Perron

We say a real matrix M is *non-negative* if all its entries are non-negative. We write $M \geq N$ if $M - N$ is non-negative. We say M is *positive* if all its entries are positive. If M is a real matrix of any order, then we define $|M|$ to be the matrix we get by replacing each entry by its absolute value.

7.12.1 Lemma. *Let A be an $n \times n$ matrix with spectral radius ρ , and suppose A is real and all its entries are positive. Suppose that θ is an eigenvalue such that $|\theta| = \rho$ and let x be an eigenvector with eigenvalue θ . Then $|x|$ is an eigenvector for A with eigenvalue ρ .*

Proof. We have

$$\rho|x| = |\theta x| = |Ax| \leq |A||x|$$

and therefore

$$A|x| \geq \rho|x|.$$

First, suppose there is a non-negative non-zero vector z such that $Az \geq \sigma z$ and $\sigma > \rho$. Then

$$A^n z \geq \sigma^n z$$

and therefore

$$\|A^n\| \geq \sigma^n$$

for all n . This implies that the spectral radius of A is at least σ , which contradicts the fact that the spectral radius equals ρ .

Now suppose that z is a non-negative non-zero vector such that $Az \geq \rho z$ and, for some index k , we have

$$e_k^T Az > \rho e_k^T z.$$

Consider the vector $z + te_k$, where t is small. Then

$$A(z + te_k) \geq \rho z + tAe_k.$$

Since all entries of A are positive, it follows that, if $i \neq k$, then

$$e_i^T A(z + te_k) > \rho e_i^T z = \rho e_i^T (z + te_k).$$

On the other hand

$$\begin{aligned} e_k^T A(z + te_k) &= e_k^T Az + te_k^T Ae_k > \rho e_k^T z + te_k^T Ae_k \\ &= \rho e_k^T (z + te_k) + t(A_{k,k} - 1). \end{aligned}$$

It follows that there are positive values of t such that

$$e_k^T A(z + te_k) > \rho e_k^T (z + te_k)$$

and, for these values of t , we have

$$A(z + te_k) > \rho z + te_k.$$

Since this is impossible, we are forced to conclude that $Az = \rho z$. □

7.12.2 Theorem. *Let A be a real square matrix with positive entries. Then the spectral radius of A is an eigenvalue of A with algebraic multiplicity 1, and corresponding eigenspace is spanned by an eigenvector with all entries positive. If θ is an eigenvalue of A not equal to ρ , then $|\theta| < \rho$.*

Proof. We have seen that there is an eigenvector x with eigenvalue ρ and all its entries non-negative. We show that the entries of any non-negative eigenvector with eigenvalue ρ must all be positive. Suppose $\rho y = Ay$ and $y \geq 0$. Then

$$\rho e_i^T y = e_i^T Ay = \sum_j A_{i,j} y_j.$$

However all entries of A are positive and y is non-negative and not zero, so the above sum is positive. As $\rho > 0$, it follows that $e_i^T y > 0$.

Next we show that ρ has geometric multiplicity 1. Assume $Ay = \rho y$, where y is not a scalar multiple of x . Then there is a real number t such that $x + ty \geq 0$ and some entry of $x + ty$ equals 0. But $x + ty$ is an eigenvector for A with eigenvalue ρ , and so we have a contradiction. We conclude that ρ has geometric multiplicity 1.

Finally we show that ρ has algebraic multiplicity 1. Suppose that $(A - \rho I)^2 w = 0$ and w is not in $\ker(A - \rho I)$. Then, replacing w by $-w$ if needed, we may assume that $x = (A - \rho I)w$ is a positive eigenvector for A with eigenvalue ρ . Note now that A^T is a positive matrix with spectral radius ρ . (It has the same minimal polynomial as A , hence has the same eigenvalues.) Let y be a positive eigenvector for A^T with eigenvalue ρ . Then $y^T(A - \rho I) = 0$, and consequently

$$y^T x = y^T (A - \rho I)w = 0.$$

But y and x are positive, and therefore $y^T x > 0$. Thus we conclude that, if $(A - \rho I)^2 w = 0$ then $w = 0$. Therefore the algebraic multiplicity of ρ is 1.

Now suppose that θ is an eigenvalue of A distinct from ρ , and let x be an eigenvector for θ . Then, using the triangle inequality,

$$|\theta| |x_i| = |(Ax)_i| = \left| \sum_j A_{i,j} x_j \right| \leq \sum_j |A_{i,j} x_j| = (A|x|)_i.$$

This implies that $|\theta| \leq \rho$. If equality holds, then

$$\left| \sum_j A_{i,j} x_j \right| \leq \sum_j |A_{i,j} x_j|.$$

Thus we have n possibly complex numbers $z_j := A_{i,j} x_j$ such that

$$\left| \sum_j z_j \right| \leq \sum_j |z_j|,$$

which implies that there is a root of unity ξ such that ξz_j is real and positive for all j . Therefore ξx is a positive eigenvector and $\theta = \rho$. \square

If $y^T A = \rho y^T$ and $Ax = \theta x$, where $\theta \neq \rho$, then $y^T x = 0$. This implies that any non-negative eigenvector for A must be an eigenvector for ρ .

7.12.3 Lemma. Let A be a real square matrix with all entries positive, and let x be a positive eigenvector for A with eigenvalue ρ , such that $\mathbf{1}^T x = 1$. If u is a non-zero non-negative vector, then

$$\lim_{n \rightarrow \infty} \frac{A^n u}{\mathbf{1}^T A^n u} = x.$$

Proof. Let x be a positive eigenvector for A with eigenvalue ρ , and let y be a positive eigenvector for A^T with eigenvalue ρ . Let B be defined by

$$B := A - \frac{\rho}{y^T x} x y^T.$$

If $Az = \theta z$ and $\theta \neq \rho$, then $y^T z = 0$ and $Bz = \theta z$. Also $Bx = 0$ and therefore if θ is an eigenvalue of B , then $|\theta| < \rho$. Consequently $\rho^{-1}B$ is a contraction. Let E be given by

$$E := \frac{1}{y^T x} x y^T.$$

Then $E^2 = E$ and $AE = EA$ and $BE = EB = 0$. Accordingly

$$(B + \rho E)^n = B^n + \rho^n E$$

and, for any vector u ,

$$A^n u - \rho^n E u = B^n u.$$

Therefore, since $\rho^{-1}B$ is a contraction, $\rho^{-n}B^n u \rightarrow 0$ as $n \rightarrow \infty$ and, provided $y^T u \neq 0$,

$$\lim_{n \rightarrow \infty} \frac{A^n u}{\mathbf{1}^T A^n u} = \lim_{n \rightarrow \infty} \frac{\rho^{-n} A^n u}{\rho^{-n} \mathbf{1}^T A^n u} = \frac{1}{\mathbf{1}^T E u} E u = \frac{1}{\mathbf{1}^T x} x.$$

- (1) Let A be a positive square matrix. Show that there is a non-negative vector x such $(I - A)x$ is non-negative and not zero if and only if A is a contraction.

Chapter 8

Geometry

We study some geometric questions.

8.1 Semilinear Forms

A *semilinear form* on a vector space V is a map from $V \times V$ to the underlying field. It maps the pair (x, y) to $\langle x, y \rangle$, and satisfies the following:

- (a) For each vector a , the map $x \mapsto \langle a, x \rangle$ is linear.
- (b) For each vector b , the map $x \mapsto \langle x, b \rangle$ is semilinear.

It follows that for all vectors x and y and all scalars a ,

$$\langle ax, y \rangle = a^\sigma \langle x, y \rangle.$$

The standard inner product on \mathbb{C}^d is semilinear; in this case σ is complex conjugation. For a wider class of examples, take a square matrix A and define

$$\langle x, y \rangle = (x^\sigma)^T Ay.$$

(For a matrix or vector M , we use M^σ to denote the result of applying σ to each entry of M .)

Since the map $\psi_a : x \mapsto \langle a, x \rangle$ is a linear map from V to the 1-dimensional space \mathbb{F} we see that either ψ_a is onto and its kernel has codimension 1 in V , or ψ_a is the zero map and its kernel is V . We denote the kernel of ψ_a by x^\perp . The radical of V (relative to our form) is the set of vectors a such that ψ_a is the zero

map. It is a subspace of V . We say that the form is *non-degenerate* if its radical is zero. The radical of an inner product is zero.

If $U \leq V$, we define

$$U^\perp = \bigcap_{u \in U} u^\perp.$$

This is again a subspace of V .

8.1.1 Lemma. *If $U \leq V$ and our form is non-degenerate, then $\dim(U) + \dim(U^\perp) = \dim(V)$.*

Let u_1, \dots, u_k be a basis for U and define a map $\rho : U \rightarrow \mathbb{F}^k$ by

$$\rho(x) = (\langle u_1, x \rangle \quad \dots \quad \langle u_k, x \rangle).$$

We see that ρ is linear and that $\ker(\rho) = U^\perp$. If ρ is not surjective, there are scalars a_1, \dots, a_k such that

$$0 = \sum_{r=1}^k a_r \langle u_r, x \rangle = \langle \sum_{r=1}^k a_r u_r, x \rangle.$$

Since our form is non-degenerate, it follows that $\sum_{r=1}^k a_r u_r = 0$ and, since u_1, \dots, u_k is a basis, $a_r = 0$ for all r . We conclude that ρ is surjective, and the lemma follows from the rank-nullity theorem. \square

We say that a subspace U of V is *isotropic* if $U \leq U^\perp$. The zero subspace is the only isotropic subspace of an inner product space.

8.2 The Classification of Forms

There are three classes of semilinear forms.

For the first, the associated automorphism is non-trivial, and

$$\langle y, x \rangle = (\langle x, y \rangle)^\sigma.$$

In this case we have a *Hermitian form*. For a Hermitian form there is a matrix H such that $(H^\sigma)^T = H$ and

$$\langle x, y \rangle = (x^\sigma)^T H y.$$

Otherwise σ is trivial. The next possibility is that

$$\langle y, x \rangle = \langle x, y \rangle.$$

In this case we have a *symmetric form*, for which there is always a symmetric matrix A such that $\langle x, y \rangle = x^T A y$. Finally we may have an *alternating form*, where

$$\langle x, x \rangle = 0$$

for all x . Here

$$0 = \langle x + y, x + y \rangle = \langle x, x \rangle + \langle x, y \rangle + \langle y, x \rangle + \langle y, y \rangle$$

and since $\langle x, x \rangle = \langle y, y \rangle$, it follows that

$$\langle y, x \rangle = -\langle x, y \rangle$$

For an alternating form there is a matrix S such that $S^T = -S$ and all diagonal entries are zero; then $\langle x, y \rangle = x^T S y$.

Alternating forms are also known as *symplectic forms*. In odd characteristic it is reasonable to describe the matrix S as skew symmetric. In even characteristic, S is symmetric with zero diagonal.

Under natural geometric assumptions it can be shown (with some effort) that the above three families of semilinear forms are the only interesting possibilities.

We say two forms \langle, \rangle_1 and \langle, \rangle_2 are *equivalent* if there is an invertible matrix M such that

$$\langle x, y \rangle_2 = \langle Mx, My \rangle_1.$$

This raises the problem of determining the equivalence classes of forms of a given type on vector space.

Over finite fields it can be shown that there is only one class of non-degenerate Hermitian forms, and only one class of non-degenerate alternating forms. It cannot be shown that there is only one class of non-degenerate symmetric forms—because this is false.

8.3 Gram Matrices

It is very useful to be able to compute the vector in U closest to a given vector x . In the previous sections, we have seen how to do this if we are given an orthogonal basis for U , or if the inner product is the dot product. We now develop the tools to solve this problem in general.

The *Gram matrix* G of a subset x_1, \dots, x_n of U is the matrix with entries given by

$$G_{i,j} = \langle x_i, x_j \rangle.$$

If $a^T = (a_1, \dots, a_n)^T$, then

$$a^T G a = \left\langle \sum_i a_i x_i, \sum_i a_i x_i \right\rangle$$

and therefore $a^T G a > 0$ for any non-zero vector a . We say a matrix G is *positive definite* if it is self-adjoint and $a^T G a > 0$ for any non-zero vector a ; if it is self-adjoint and $a^T G a \geq 0$ for all a , then G is *positive semidefinite*. We have just seen that Gram matrices are positive semidefinite.

8.3.1 Lemma. *A set of vectors in an inner product space is linearly independent if and only if their Gram matrix is invertible.*

Proof. Suppose G is the Gram matrix for x_1, \dots, x_n . Then the entries of Ga are the inner products

$$\langle x_r, \sum_r a_r x_r \rangle$$

Hence if U is the span of the vectors x_1, \dots, x_n , then $Ga = 0$ if and only if $\sum_r a_r x_r = 0$. Thus $\ker(G)$ is zero if and only if x_1, \dots, x_n are linearly independent. \square

8.4 Equiangular Lines

We work in the vector space V , which is \mathbb{R}^d or \mathbb{C}^d with the usual Euclidean inner product. If x and y are nonzero vectors, the *cosine* of the angle between the lines spanned by x and y is

$$\frac{|\langle x, y \rangle|}{\|x\| \|y\|}.$$

We will often work with the squared cosine

$$\frac{\langle x, y \rangle \langle y, x \rangle}{\langle x, x \rangle \langle y, y \rangle}.$$

A set of lines in V is *equiangular* if the cosine of the angle between any two distinct lines is the same.

8.4.1 Theorem. *The maximum size of a set of equiangular lines in \mathbb{C}^d is d^2 ; in \mathbb{R}^d it is $\binom{d+1}{2}$.*

Proof. Suppose we have lines spanned by unit vectors x_1, \dots, x_m . Define matrices P_1, \dots, P_m by

$$P_r = x_r x_r^*.$$

Then P_r represents orthogonal projection onto the line spanned by x_r , and if $r \neq s$,

$$\langle P_r, P_s \rangle = \text{tr}(P_r P_s) = \langle x_r, x_s, \langle \cdot, x_s \rangle x_r \rangle = |\langle x_r, x_s \rangle|^2.$$

We assume that $\alpha = |\langle P_r, P_s \rangle|$. We see also that $\langle P_r, P_r \rangle = 1$ for all r .

The projections P_r lie in the space of Hermitian matrices. If G is their Gram matrix, then

$$G = (1 - \alpha^2)I + \alpha^2 J.$$

We can prove, in a number of ways, that G is invertible, which implies that the matrices P_1, \dots, P_m form a linearly independent set in the space of Hermitian matrices. We complete the proof by noting that this space has dimension d^2 (over \mathbb{C}) and the dimension in the real case is $\binom{d+1}{2}$. \square

In \mathbb{R}^2 it is easy to find three lines with pairwise cosine $1/2$, and the diagonals of the icosahedron give six lines with pairwise cosine $1/\sqrt{5}$. Examples of sets of size $\binom{d+1}{2}$ are known in \mathbb{R}^d when $d = 7$ and $d = 23$. In the complex case, examples of tight sets are known for d in $\{1, \dots, 15, 19, 24, 35, 48\}$.

8.5 Tight Frames

Suppose we have a set of equiangular lines of maximum size. Then the associated projections P_1, \dots, P_m form a basis for the space of Hermitian matrices. Hence there are scalars c_r such that

$$I = \sum_r c_r P_r.$$

If we multiply both sides by P_k and take traces, we get

$$1 = (1 - \alpha^2)c_k + \alpha^2 \sum_r c_r.$$

It follows that $c_1 = \dots = c_m$ and hence that

$$I = \frac{d}{m} \sum_r P_r.$$

In a slightly different format, we have established that if x_1, \dots, x_m are unit vectors spanning a set of equiangular lines of maximum size, then

$$\sum_r x_r x_r^* = \frac{m}{d} I.$$

Such a set of vectors is an example of a *tight frame*.

We will see that tight frames are more common than set of lines meeting the absolute bound. Consider a set of projections P_1, \dots, P_m corresponding to a set of equiangular lines with squared cosine α^2 , and define

$$M = \sum_r P_r - \frac{m}{d} I.$$

Then

$$\begin{aligned} 0 \leq \langle M, M \rangle &= \langle \sum_r P_r, \sum_r P_r \rangle - \frac{2m}{d} \langle \sum_r P_r, I \rangle + \frac{d^2}{m^2} \text{tr}(I) \\ &= m + m(m-1)\alpha^2 - \frac{m^2}{d}. \end{aligned}$$

If equality holds we have

$$\alpha^2 = \frac{m-d}{md-d}.$$

This yields the following, sometimes known as the *relative bound*.

8.5.1 Theorem. *If there is a set of m lines in \mathbb{F}^d with squared cosine α^2 , where $d\alpha^2 < 1$, then*

$$m \leq \frac{d - d\alpha^2}{1 - d\alpha^2}.$$

If equality holds, then a set of unit vectors spanning the lines forms a tight frame. \square

Note that if we have d^2 lines in \mathbb{C}^d , then $\alpha^2 = (d+1)^{-1}$, and for $\binom{d+1}{2}$ lines in \mathbb{R}^d , then $\alpha^2 = (d+2)^{-1}$.

8.6 Another Gram Matrix

Suppose x_1, \dots, x_m form a tight frame in dimension d . Then

$$\sum_r x_r x_r^* = \frac{m}{d} I.$$

If U is the $d \times m$ matrix with the vectors x_1, \dots, x_m as its columns then we have

$$UU^* = \sum_r x_r x_r^* = \frac{m}{d} I,$$

which implies that the rows of U are orthogonal (and of the same length).

Set $H = U^*U$. Then

$$H^2 = U^*UU^*U = \frac{m}{d}U^*U = \frac{m}{d}H$$

and therefore the minimal polynomial of H divides

$$t\left(t - \frac{m}{d}\right).$$

(If the minimal polynomial is a proper divisor of this polynomial that $H = 0$ or $H = I$.) We can write H as $I + \alpha S$, where S is Hermitian with diagonal entries zero and all off-diagonal entries have absolute value 1. (In the real case, this means the off-diagonal entries are ± 1 .) The eigenvalues of S are

$$\frac{1}{\alpha} \left(\frac{m}{d} - 1 \right), \quad -\frac{1}{\alpha}$$

with respective multiplicities d and $m - d$.

8.7 The Orthogonal Group

Let V be a vector space with a bilinear form. We say that an endomorphism A of V preserves the form if $\langle Ax, Ay \rangle = \langle x, y \rangle$, for all x and y . If the form is symmetric and the characteristic of our field is odd, then

$$\langle x, y \rangle = \frac{1}{2}(\langle x + y, x + y \rangle - \langle x, x \rangle - \langle y, y \rangle).$$

Hence A preserves the form if and only if $\langle Ax, Ax \rangle = \langle x, x \rangle$ for all x .

Now assume V is \mathbb{R}^n and that our form is the dot product. A matrix which preserves dot product is called *orthogonal*. If v and w are orthogonal vectors in V and A is orthogonal, then Av and Aw are orthogonal.

8.7.1 Lemma. *A matrix A is orthogonal if and only if $A^T A = I$.*

Proof. If v_1, \dots, v_n is an orthogonal basis for V , then so is Av_1, \dots, Av_n . Since the standard basis e_1, \dots, e_n for V is orthogonal, it follows that Ae_1, \dots, Ae_n is an orthogonal set of vectors. Therefore the columns of an orthogonal matrix A form an orthogonal basis. This also implies that

$$A^T A = I.$$

Since A is square, we see that $A^T = A^{-1}$ and $AA^T = I$. Conversely, if $A^T = A^{-1}$, then

$$\langle Av, Aw \rangle = (Av)^T Aw = v^T A^T Aw = v^T w = \langle v, w \rangle. \quad \square$$

We see from this result that, if A is orthogonal, then its columns form an orthonormal set. Also, if A is orthogonal, then $A^T = A^{-1}$ and therefore $AA^T = I$. Hence the rows of A also form an orthonormal set.

We consider the complex version of orthogonal matrices. A complex matrix is *unitary* if it preserves the complex dot-product. This means that

$$y^* x = (Ay)^* (Ax) = y^* A^* Ax$$

for all x and y , and hence that

$$A^* A = I.$$

A real matrix is unitary if and only if it is orthogonal.

We turn to examples of orthogonal matrices. Any permutation matrix is orthogonal, and a diagonal matrix A is orthogonal if and only if $A_{i,i} = \pm 1$ for all i . The matrices

$$\begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix}$$

are orthogonal, for any value of θ . It is easy to verify that the product of two orthogonal matrices is orthogonal, and that the inverse of an orthogonal matrix is orthogonal. Therefore the set of all orthogonal matrices is an example of a *group*, known as the *orthogonal group*.

8.8 Skew-Symmetric to Orthogonal

We define a matrix A to be *skew symmetric* if $A^T = -A$ and $A_{i,i} = 0$ for all i . (The last condition is only needed if our field has characteristic two.) The set of $n \times n$ skew-symmetric matrices is a subspace of the space of square matrices.

8.8.1 Lemma. *If S is a real skew-symmetric matrix, then $(I - S)^{-1}(I + S)$ is orthogonal.*

Proof. We first show that $I - S$ is invertible for all real t . Suppose $x \neq 0$ and $Ax = \theta x$. Then

$$\theta x^T x = x^T Ax = (A^T x)x = (-Ax)^T x = (-\theta x)^T x = -\theta x^T x.$$

It follows that 0 is the only possible real eigenvalue for S . Therefore $I - S$ is invertible for all real t and we can define

$$M := (I - S)^{-1}(I + S).$$

The matrices $I + S$ and $I - S$ commute, and from this it follows that $I + S$ and $(I - S)^{-1}$ commute. Hence we find that

$$\begin{aligned} M^T &= (I + S^T)(I - S^T)^{-1} = (I - S)(I + S)^{-1} \\ &= (I + S)^{-1}(I - S) \\ &= M^{-1}. \end{aligned}$$

Therefore M is orthogonal. □

The matrix M above is sometimes known as the *Cayley transform* of S . Note that, since tS is skew-symmetric if S is, the matrix

$$(I - tS)^{-1}(I + tS)$$

is orthogonal for real t .

(1) If H is hermitian and $S = iH$, show that $(I - S)^{-1}(I + S)$ is unitary.

8.9 Reflections

Suppose a is a fixed non-zero vector in V . Define the map ρ_a by

$$\rho_a(v) = v - 2 \frac{\langle a, v \rangle}{\langle a, a \rangle} a.$$

Note that ρ_a is the sum of two linear mappings (the identity and a scalar multiple of the orthogonal projection onto the line spanned by a) and therefore it is linear. We check that

$$\rho_a(a) = -a$$

and, using this, that

$$\rho_a^2 = I.$$

If $v \in a^\perp$, then $\rho_a(v) = v$. It follows that ρ_a corresponds to the geometric operation of reflection in the hyperplane perpendicular to a .

We have

$$\langle \rho_a(v), \rho_a(v) \rangle = \langle v, v \rangle - 4 \frac{\langle a, v^2 \rangle}{i p a a} + 4 \frac{\langle a, v \rangle^2}{\langle a, a \rangle} \langle a, a \rangle = \langle v, v \rangle.$$

Therefore ρ_a is orthogonal. The matrix R_a representing it is given by

$$R_a = I - \frac{2}{\langle a, a \rangle} a a^T.$$

If v and w have the same length then

$$\langle v - w, v + w \rangle = 0.$$

Therefore R_{v-w} fixes $v + w$ and maps $v - w$ to $w - v$. Consequently

$$R_{v-w}(2v) = R_{v-w}((v + w) + (v - w)) = 2w,$$

and so, after a very modest amount of extra work, we find that R_{v-w} swaps v and w .

8.9.1 Theorem. *Every non-identity orthogonal matrix is a product of at most n matrices R_a .*

Proof. If A is a matrix, let $F(A)$ be the subspace

$$\{v \in V : Av = v\}.$$

We prove by induction that A is the product of at most $\dim(V) - \dim(F(A))$ matrices R_a .

Suppose A is orthogonal and $\dim(F(A)) = k$. If $k = \dim V$, then $A = I$. Suppose $k < \dim V$, and let v be a vector in V such that $Av \neq v$. If $w := Av$ and $Ax = x$, then $\langle v, x \rangle = \langle w, x \rangle$ and so $\langle v - w, x \rangle = 0$. Therefore $F(A) \subseteq (v - w)^\perp$, and R_{v-w} fixes each vector in $F(A)$. Now R_{v-w} swaps v and w , whence the product $R_{v-w}A$ fixes each vector in $F(A)$, and fixes v . As $v \notin F(A)$, we see that

$$\dim(F(R_{v-w}A)) > \dim(F(A)).$$

The lemma follows. □

A matrix A is an involution if $A^2 = I$. Diagonal matrices with diagonal entries equal to ± 1 provide a fairly trivial class of examples. If P is an idempotent then

$$(I - 2P)^2 = I - 4P + 4P = I,$$

and thus $I - 2P$ is an involution.

8.9.2 Theorem. *Every orthogonal matrix is the product of two involutions.*

Proof. We actually prove a stronger result: A and A^{-1} are similar if and only if A is the product of two involutions. Since any square matrix is similar to its transpose, orthogonal matrices satisfy this condition.

Suppose $S^2 = T^2 = I$ and $A = ST$. Then $(ST)(TS) = I$ and

$$S^{-1}AS = SAS = S(ST)S = TS.$$

Therefore a product of two involutions is similar to its inverse.

So assume now that A and A^{-1} are similar and let F be the Frobenius normal form of A . By ??, there is a permutation matrix T such that $T^2 = I$ and

$$F^{-1} = TFT.$$

Then $I = FTFT$, whence FT and T are involutions whose product is F . As any matrix that is similar to an involution is an involution, the general result follows. \square

Chapter 9

Positive Semidefinite Matrices

9.1 Factorizing Positive Semidefinite Matrices

If H is a matrix with linearly independent columns, then the product $H^T H$ is the Gram matrix for a basis of $\text{col}(H)$ and therefore it is positive definite. Our next result provides a converse to this.

9.1.1 Theorem. *If G is a positive definite matrix, there is a lower triangular matrix L with diagonal entries equal to 1 and a diagonal matrix D with positive diagonal entries, such that $LGL^T = D$.*

Proof. If G is positive definite, then $e_i^T G e_i > 0$ for all i ; hence the diagonal entries of G are positive.

Since L and G are invertible, $D = LGL^T$ is necessarily invertible. We must show that L exists. We write G in partitioned form:

$$G = \begin{pmatrix} a & b^T \\ b & G_1 \end{pmatrix}.$$

If we also define

$$L_1 = \begin{pmatrix} 1 & 0 \\ -a^{-1}b & I \end{pmatrix}$$

then

$$L_1 G L_1^T = \begin{pmatrix} a & 0 \\ 0 & G_1 - a^{-1} b b^T \end{pmatrix}.$$

Note that $a \neq 0$, because G is positive definite. It follows from the exercises below that $G_1 - a^{-1} b b^T$ is positive definite. By induction, we have that there is

a lower triangular matrix L_2 with diagonal entries equal to 1 such that

$$L_2(G - a^{-1}bb^T)L_2^T$$

is diagonal. Taking L to be given by

$$L := \begin{pmatrix} 1 & 0 \\ 0 & L_2 \end{pmatrix} L_1,$$

our result follows. \square

This result implies that $G = L^{-1}DL^{-T}$. Since the diagonal entries of D are positive, there is a unique non-negative diagonal matrix $D^{1/2}$ such that $(D^{1/2})^2 = D$ and therefore

$$G = (L^{-1}D^{1/2})(L^{-1}D^{1/2})^T.$$

A factorization of a positive-definite matrix G in the form MM^T , where M is lower triangular with positive diagonal entries, is known as a *Cholesky factorization*. Any reasonable software package for linear algebra will have a command to compute the matrix M from G .

If G is presented as a matrix $X^T X$ and $LGL^T = D$, then

$$(XL^T)^T(XL^T) = D,$$

whence we see that the columns of XL^T are orthogonal (with respect to the dot product). Thus they form an orthogonal basis for $\text{col}(X)$, and so we may use the Cholesky decomposition to find orthogonal bases. We illustrate this in the next section.

We record an important property of positive definite matrices—it is basically a reformulation of the definition.

9.1.2 Lemma. *If A is a positive definite matrix, the bilinear form*

$$\langle x, y \rangle = x^T Ay$$

is an inner product.

Proof. Exercise. \square

- (1) If G is positive definite and the columns of L are linearly independent, show that LGL^T is positive definite.
- (2) Show that a principal submatrix of a positive definite matrix is positive definite.
- (3) Prove that if G has Cholesky factorizations MM^T and NN^T , then $M = N$.

9.2 Computing Cholesky

The Cholesky decomposition of a positive definite matrix can be useful, in particular it may be used to find orthogonal bases. In this section we describe an algorithm for computing the Cholesky factorization using elementary row operations. (But outside linear algebra courses, we recommend using methods based on the QR-factorization, which we address later. Our point is that we can carry out Gram-Schmidt by using Gaussian elimination.)

As a first step, we need to note one consequence of Theorem 9.1.1. This result shows that if G is positive definite, then by successively subtracting multiples of higher rows from lower rows, we can convert G to an invertible upper triangular matrix. The product of the elementary matrices corresponding to these operations is the lower triangular matrix L . Our next result asserts that if we use elementary operations as described to bring G to row echelon form, we obtain the Cholesky factorization of G .

9.2.1 Lemma. *Let G be a positive definite matrix. If K is lower triangular with diagonal entries equal to 1 and KG is upper triangular, then $KG = DK^{-T}$, where D is a diagonal matrix with positive diagonal entries.*

Proof. Suppose that K is lower triangular with diagonal entries equal to 1, and that $KG = DM$, where D is diagonal and M is upper triangular, with diagonal entries 0 or 1. Then

$$K GK^T = DMK^T.$$

Here the left side is a symmetric matrix, while the right side is the product of three upper triangular matrices, and is therefore upper triangular. It follows that MK^T is diagonal. Since $K GK^T$ is invertible, both D and MK^T are invertible. Therefore $MK^T = I$. Finally $K GK^T$ is positive definite and equal to D . So D is positive definite, and therefore its diagonal entries are positive. \square

Suppose we are given a Gram matrix G . If we bring the partitioned matrix

$$\begin{pmatrix} G & I \end{pmatrix}$$

to row-echelon form, then the resulting matrix equals

$$\begin{pmatrix} LG & L \end{pmatrix}.$$

As noted at the end of the previous section, if $G = X^T X$, then the columns of XL^T are orthogonal (with respect to the dot product). The i -th column of

XL^T is a linear combination of the first i columns of X and consequently the columns of XL^T are the orthogonal set we would compute using the usual approach to Gram-Schmidt. (Using exact arithmetic—in fact we have developed the so-called modified Gram-Schmidt method.)

We first illustrate this in \mathbb{R}^n , with the dot product. The row echelon form of the partitioned matrix

$$M = \begin{pmatrix} X^T X & X^T \end{pmatrix}$$

is

$$\begin{pmatrix} LX^T X & LX^T \end{pmatrix}$$

and so the transposes of the rows of LX^T are an orthogonal basis for the column space of X . Suppose for example that

$$x_1 = \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}, \quad x_2 = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}, \quad x_3 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}.$$

Let X be the matrix with x_1, x_2 and x_3 as its columns. Then

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

has row echelon form

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

Hence

$$XL^T = \begin{pmatrix} 1 & -\frac{1}{2} & \frac{2}{3} \\ 1 & \frac{1}{2} & -\frac{2}{3} \\ 0 & 1 & \frac{2}{3} \end{pmatrix}$$

and its columns are an orthogonal basis for $\text{col}(X)$.

9.3 Polynomial Examples

We consider the situation where we want to find an orthogonal basis for an inner product space of polynomials. By way of example, we take V to be the space of all polynomials, with inner product:

$$\langle p, q \rangle := \int_0^\infty p(x)q(x)e^{-x} dx.$$

Let U be the subspace consisting of the polynomials with degree at most n , let p_0, \dots, p_n be basis for U and let G be the Gram matrix of this basis. (Thus the rows and columns of G are indexed by $0, 1, \dots, n$, rather than $1, \dots, n$ —good news for C programmers anyway.)

If $[q]$ denotes the coordinate vector of q in U relative to the given basis, then

$$[p]^T G[q] = \langle p, q \rangle.$$

Suppose $LGL^T = D$. Then

$$e_i^T LGL^T D = e_i^T D e_j$$

whence the columns of L^T are the coordinate vectors of an orthogonal basis for U .

Turning to a concrete case, suppose U is the space of polynomials with degree at most three. We start with the basis $1, x, x^2, x^3$. It can be shown (by integration by parts) that

$$\int_0^\infty x^n e^{-x} dx = n!,$$

and therefore the Gram matrix of this set of polynomials is

$$G = \begin{pmatrix} 1 & 1 & 2 & 6 \\ 1 & 2 & 6 & 24 \\ 2 & 6 & 24 & 120 \\ 6 & 24 & 120 & 720 \end{pmatrix}.$$

Let M be given by

$$M = \begin{pmatrix} 1 & 1 & 2 & 6 & 1 & 0 & 0 & 0 \\ 1 & 2 & 6 & 24 & 0 & 1 & 0 & 0 \\ 2 & 6 & 24 & 120 & 0 & 0 & 1 & 0 \\ 6 & 24 & 120 & 720 & 0 & 0 & 0 & 1 \end{pmatrix}.$$

We convert the first four columns to an upper triangular matrix:

$$\begin{pmatrix} 1 & 1 & 2 & 6 & 1 & 0 & 0 & 0 \\ 0 & 1 & 4 & 18 & -1 & 1 & 0 & 0 \\ 0 & 0 & 4 & 36 & 2 & -4 & 1 & 0 \\ 0 & 0 & 0 & 36 & -6 & 18 & -9 & 1 \end{pmatrix},$$

and thus obtain the following set of four orthogonal polynomials:

$$1, \quad x-1, \quad x^2-4x+2, \quad x^3-9x^2+18x-6.$$

9.4 Positive Semidefinite Matrices

We develop some further properties of positive semidefinite matrices.

9.4.1 Lemma. *If A and B are positive semidefinite, so is $A + B$. If A is positive and B is positive definite, then $A + B$ is positive definite. \square*

We leave the proof as an exercise. Note that it implies that if A is positive semidefinite, then $A + I$ is positive definite.

9.4.2 Lemma. *A self-adjoint matrix is positive semidefinite if and only if its eigenvalues are non-negative. It is positive definite if and only if its eigenvalues are positive.*

Proof. If x is an eigenvector of A with eigenvalue θ , then $x^T A x = \theta x^T x$, and therefore if A is positive semidefinite, its eigenvalues are non-negative. If A is positive definite then 0 is not an eigenvalue.

Suppose we have the spectral decomposition

$$A = \sum_{\theta} \theta E_{\theta}.$$

Each projection E_{θ} is positive semidefinite, because

$$x^T E_{\theta} x = x^T E_{\theta}^2 x = x^T E_{\theta}^T E_{\theta} x = \|E_{\theta} x\|^2.$$

If each eigenvalue of A is non-negative, it follows that $x^T A x$ is a sum of non-negative terms $\theta x^T E_{\theta} x$, and therefore $x^T A x \geq 0$.

If the eigenvalues of A are positive, we see that $x^T A x = 0$ if and only if $x^T E_{\theta} x = 0$ for each eigenvalue θ . Hence

$$0 = \sum_{\theta} x^T E_{\theta} x = x^T \left(\sum_{\theta} E_{\theta} \right) x = x^T I x,$$

and therefore $x = 0$. Consequently A is positive definite. \square

Note that I_n has 2^n distinct square roots, that is, there are 2^n matrices S such that $S^2 = I$. However it has only one positive semidefinite square root. This is typical:

9.4.3 Corollary. *If A is positive semidefinite, there is a unique positive semidefinite matrix S such that $S^2 = A$.*

Proof. Using the spectral decomposition we have

$$A = \sum_{\theta} \theta E_{\theta},$$

where the sum is over all eigenvalues of A . If A is positive semidefinite, its eigenvalues are non-negative and we may define S by

$$S = \sum_{\theta} \sqrt{\theta} E_{\theta}.$$

Since the eigenvalues of S are non-negative, it is positive semidefinite.

We turn to uniqueness. Let T be a positive semidefinite square root of A and suppose x is an eigenvector for A . If $Ax = 0$ then $T^2x = 0$, so $x^T TTx = 0$ and therefore $Tx = 0$. Assume now that $Ax = \sigma^2 x$, where $\sigma > 0$, then

$$0 = (T^2 - \sigma^2 I)x = (T - \sigma I)(T + \sigma I)x.$$

If the subspace spanned by x is T -invariant, it follows that $Tx = \pm\sigma x$ and x is an eigenvector for T . Otherwise x and Tx span a T -invariant subspace on which T acts with minimal polynomial $t^2 - \sigma^2$. If $(T - \sigma I)x \neq 0$ then $y = (T - \sigma I)x$ is an eigenvector for $T + \sigma I$ with eigenvalue $-\sigma$. Therefore if T is positive semidefinite and $Ax = \sigma^2 x$, then $Tx = \sigma x$.

Thus we have shown that, if $Ax = \sigma^2 x$ then $Tx = \sigma x$. Since the eigenvectors of A span, this shows that T is determined by A . \square

The next result is known as the *polar decomposition* of a matrix. It is analogous to the fact that each complex number is the product of a positive real number and a complex number with norm 1.

9.4.4 Theorem. *If A is a square matrix, there is a positive semidefinite matrix M and an orthogonal matrix Q such that $A = MQ$.*

Proof. We use the singular value decomposition, which yields that

$$A = Y \Sigma X^T,$$

where X and Y are orthogonal and Σ is positive semidefinite. Hence

$$A = Y \Sigma Y^T Y X^T,$$

where $Y \Sigma Y^T$ is positive semidefinite and $Y X^T$ is orthogonal. \square

Note that $AA^T = (MQ)(MQ)^T = M^2$; hence the positive definite factor in the above theorem is unique, and the orthogonal factor is unique if A is invertible.

Chapter 10

Tensors

10.1 Tensor Products

The *tensor product* $U \otimes V$ of two vector spaces U and V over \mathbb{F} is defined as a quotient space. We start with the space of all finitely supported functions $\mathbb{F}^{U \times V}$, the tensor product is the quotient of this subspace modulo the subspace spanned \mathcal{R} by vectors of the following forms:

- (a) $a(u, v) - (au, v), a(u, v) - (u, av)$ for $a \in \mathbb{F}$ and (u, v) in $U \times V$.
- (b) $(u_1 + u_2, v) - (u_1, v) - (u_2, v)$ for $u_1, u_2 \in U, v \in V$.
- (c) $(u, v_1 + v_2) - (u, v_1) - (u, v_2)$ for $u \in U, v_1, v_2 \in V$.

(Here we are using formal sum of finitely many terms to represent elements of $\mathbb{F}^{U \otimes V}$.) We denote the image of (u, v) in $U \otimes V$ by $u \otimes v$. The map that sends (u, v) to $u \otimes v$ is bilinear.

For finite-dimensional vector spaces, there is no harm in identifying the tensor product with Kronecker product.

The tensor product is not commutative, the spaces $U \otimes V$ and $V \otimes U$ are isomorphic but not equal. The tensor product is associative, in that

$$(U \otimes V) \otimes W \cong U \otimes (V \otimes W).$$

The vectors of the form $u \otimes v$ are known as *pure tensors*; they span $U \otimes V$ but do not form a basis. We note that a scalar times a pure tensor is a pure tensor, and so any element of $U \otimes V$ can be expressed as a sum of pure tensors. If $\alpha \in U \otimes V$, we define the *tensor rank* of α to be the least integer r such that α

can be expressed as the sum of r pure tensors, that is, the least integer r such that

$$\alpha = \sum_{i=1}^r u_i \otimes v_i.$$

The key property of the tensor product is that it allows us to deal with linear maps in place of multilinear maps (at the cost of increasing dimensions). Thus if we have a bilinear map

$$\beta: U \times V \rightarrow W,$$

then there is a linear map $\hat{\beta}$ from $U \otimes V$ to W such that

$$\hat{\beta}(u \otimes v) = \beta(u, v).$$

If A and B are linear maps defined on U and V respectively, we define their tensor product $A \otimes B$ by

$$(A \otimes B)(u \otimes v) = Au \otimes Bv.$$

If we have inner products defined on U and V , we can define

$$\langle (u_1 \otimes v_1), (u_2 \otimes v_2) \rangle = \langle u_1, u_2 \rangle \langle v_1, v_2 \rangle.$$

This is a consequence of our definition of the tensor products of maps, because the maps $\langle u_1, ? \rangle$ and $\langle u_2, ? \rangle$ are elements of U^* .

The field \mathbb{F} is a 1-dimensional vector space and so the tensor product $\mathbb{F} \otimes V$ is defined. The map that sends $1 \otimes v$ to v is an isomorphism. If $\psi \in U^*$, then $\psi \times I$ is a linear map from $U \otimes V$ to $\mathbb{F} \otimes V$, and hence it determines a linear map from $U \otimes V$ to V . We will usually identify these two maps.

10.2 Quadratic Tensors

We investigate properties of elements of the tensor product $U \otimes V$.

10.2.1 Lemma. *If $\alpha \in U \otimes V$ has tensor rank r and $\alpha = \sum_{i=1}^r u_i \otimes v_i$ for some vectors u_1, \dots, u_r and v_1, \dots, v_r , then both of these sets of vectors are linearly independent.*

We leave the proof of this as an exercise.

10.2.2 Lemma. *If $\alpha \in U \otimes V$ has tensor rank r and*

$$\alpha = \sum_{i=1}^r u_i \otimes v_i = \sum_{i=1}^r x_i \otimes y_i$$

then

$$\text{span}\{u_1, \dots, u_r\} = \text{span}\{x_1, \dots, x_r\}, \quad \text{span}\{v_1, \dots, v_r\} = \text{span}\{y_1, \dots, y_r\}.$$

Proof. There are vectors ψ_1, \dots, ψ_r in U^* such that $\psi_i(u_j) = \delta_{i,j}$. So the image of α under the map $\psi_k \otimes 1$ is v_k , according to the first expression for α , and its image is

$$\sum_{i=1}^r \psi_i(x_i) y_i.$$

This shows that $v_k \in \text{span}\{y_1, \dots, y_r\}$, and now everything follows. \square

The previous results are analagous to properties of the usual rank of a matrix. This is no accident:

10.2.3 Theorem. *For any two vector spaces U and V , the spaces $\text{Lin}(U, V)$ and $U^* \otimes V$ are isomorphic. Under this isomorphism elements of $U^* \otimes V$ with tensor rank r map to operators with rank r .*

Proof. If $\psi \in U^*$ and $v \in V$, let the map $\lambda_{\psi,v}$ be given by

$$\lambda_{\psi,v}(u) = \psi(u)v.$$

This assigns a linear map to each pure tensor in $U^* \otimes V$ and hence gives us a linear map from $U^* \otimes V$ to $\text{Lin}(U, V)$. Denote this map by Λ .

We show that Λ is onto. The first step is to show that each linear map in $\text{Lin}(U, V)$ with rank one is the image of a pure tensor. We leave this as an exercise.

The second step is to show that any $m \times n$ matrix can be written as a sum of rank-one matrices. Suppose A is $m \times n$. If $A \neq 0$, there are vectors x and y such that $x^T A y \neq 0$, and so we may assume that we have vectors x and y such that $x^T A y = 1$. Define

$$B = A - A x y^T A^T.$$

Each column of $A x y^T A^T$ is a scalar multiple of Ax , and it follows that the column space of B is contained in the column space of A . Next, $Ax \neq 0$ but

$$Bx = Ax - A x y^T A^T x = Ax - (x^T A y) Ax = 0.$$

Therefore the column space of B is a proper subspace of the column space of A and so $\text{rk}(B) < \text{rk}(A)$. On the other hand

$$\text{rk}(A) = \text{rk}(B + Axy^T A) \leq \text{rk}(B) + \text{rk}(Axy^T A^T) \leq \text{rk}(B) + 1$$

and we conclude that $\text{rk}(B) = \text{rk}(A) - 1$. It follows by induction that A can be expressed as the sum of r rank-one matrices. \square

Note that it is immediate that a matrix with m rows is the sum of m rank-one matrices, and we can use this to provide a simple proof of the isomorphism in the above theorem. However the relation between tensor rank and the usual rank is important.

10.2.4 Theorem. *We have*

$$\dim(U \otimes V) = \dim(U) \dim(V).$$

Proof. If u_1, \dots, u_m and v_1, \dots, v_n are basis for U and V respectively, then the pure tensors $u_i \otimes v_j$ span $U \otimes V$. This shows that

$$\dim(U \otimes V) \leq \dim(U) \dim(V).$$

Proof. Suppose that the pure tensors $u_i \otimes v_j$ are linearly dependent. Then there linearly independent vectors u_1, \dots, u_r in U and vectors w_1, \dots, w_r in V such that

$$0 = \sum_{i=1}^r u_i \otimes w_i.$$

As before, choose elements f_1, \dots, f_r in U^* such that $\psi_i(u_j) = \delta_{i,j}$. If we apply $\psi_k \otimes 1$ to each side of the above expression, we get

$$0 = w_k. \quad \square$$

10.3 Cubic Tensors

Consider a tensor α in $U \otimes V \otimes W$ given by

$$\alpha = \sum_{i=1}^r u_i \otimes \beta_i$$

where $\beta_1, \dots, \beta_r \in V \otimes W$. The subspace $\mathcal{C}(\alpha)$ of $V \otimes W$ spanned by the tensors β_1, \dots, β_r is an invariant of α . Define the *order* of a subspace of $V \otimes W$ to be the

least integer s such that it is contained in the span of s pure tensors. If $\mathcal{C}(\alpha)$ has order s , then there are vectors v_1, \dots, v_s in V and pure tensors $\gamma_1, \dots, \gamma_s$ in $V \otimes W$ such that

$$\alpha = \sum_{i=1}^s x_i \otimes \gamma_i$$

Hence the tensor rank of α is at most s . Since no proper subset of $\gamma_1, \dots, \gamma_s$ spans $\mathcal{C}(\alpha)$, it follows that s is the tensor rank of α .

We give one example of the order of a subspace. Identify $V \otimes W$ with the vector space of matrices of order $\dim(V)$ times $\dim(W)$. If \mathcal{C} is the space of upper-triangular 2×2 matrices then \mathcal{C} has dimension three and order four.

For quadratic tensors, we have the following theorem.

10.3.1 Theorem. *For vector spaces V and W , the set*

$$S_k := \{T \in V \otimes W \mid rk(T) \leq k\}$$

is closed; i.e., if $\lim_{i \rightarrow \infty} T_i = T$ and $rk(T_i) \leq k$, then $rk(T) \leq k$.

Proof. Each $T \in V \otimes W$ is associated with a matrix A whose rank is equal to the tensor rank of T . Hence the sets S_k are determined by algebraic equations and are closed. \square

The set of matrix of rank at most r is a closed set, and so the limit of any sequence of matrices with rank at most r is a matrix with rank at most r . Tensor rank is in general less well behaved. Let V be \mathbb{R}^2 with the standard basis e_1, e_2 .

10.3.2 Lemma. *The element*

$$T := e_1 \otimes e_1 \otimes e_1 + e_1 \otimes e_2 \otimes e_2 + e_2 \otimes e_1 \otimes e_2$$

of $\mathbb{R}^2 \otimes \mathbb{R}^2 \otimes \mathbb{R}^2$ has tensor rank three, but is the limit of a sequence of tensors with rank at most two.

Proof. Define

$$T_\lambda := \lambda^{-1}[e_1 \otimes e_1 \otimes (-e_2 + \lambda e_1) + (e_1 + \lambda e_2) \otimes (e_1 + \lambda e_2) \otimes e_2].$$

Then

$$T_\lambda - T = \lambda e_2 \otimes e_2 \otimes e_2,$$

whence T_λ converges to T as $\lambda \rightarrow 0$.

The only difficulty is to verify that T has tensor rank three. Suppose by way of contradiction that

$$T = (\alpha_1 e_1 + \alpha_2 e_2) \otimes b \otimes c + (\mu_1 e_1 + \mu_2 e_2) \otimes v \otimes w.$$

Then

$$T = e_1 \otimes (\alpha_1 b \otimes c + \mu_1 v \otimes w) + e_2 \otimes (\alpha_2 b \otimes c + \mu_2 v \otimes w).$$

Comparing this with the definition of T , we deduce that

$$\begin{aligned} e_1 \otimes e_1 + e_2 \otimes e_2 &= \alpha_1 b \otimes c + \mu_1 v \otimes w \\ e_1 \otimes e_2 &= \alpha_2 b \otimes c + \mu_2 v \otimes w. \end{aligned}$$

The two vectors on the left in these expressions are linearly independent, and therefore these equations imply that $b \otimes c$ and $v \otimes w$ are linearly independent and that they are linear combinations of the vectors on the left

Now we use the isomorphism between $\mathbb{R}^2 \otimes \mathbb{R}^2$ and $\text{Mat}_{2 \times 2}(\mathbb{R})$. The image of the span of the vectors on the left consists of all matrices of the form

$$\begin{pmatrix} x & y \\ 0 & x \end{pmatrix}$$

All rank-one matrices of this form must have x equal to 0, and so the rank-one matrices of this form span a 1-dimensional space. It follows that $b \otimes c$ and $v \otimes w$ are linearly dependent. This the contradiction we wanted—we conclude that the tensor rank of T is three. \square

10.4 Multiplication

Let M be the space of $n \times n$ matrices over some field. Matrix multiplication defines a linear map from $V \otimes V$ to V . By Theorem 10.2.3 we have

$$\mathcal{L}(V \otimes V, V) \cong M^* \otimes M^* \otimes M,$$

and so matrix multiplication can be viewed as a particular element of this space. More concretely, if the elements $E_{i,j}$ form a basis for M and $\epsilon_{i,j}$ denotes the element of M^* that sends a matrix to its ij -entry, then

$$AB = \sum_{i,j,k} \epsilon_{i,j}(A) \epsilon_{j,k}(B) E_{i,k}$$

and so matrix multiplication corresponds to the tensor

$$\sum_{i,j,k} \epsilon_{i,j} \otimes \epsilon_{j,k} \otimes E_{i,k}.$$

This is a sum of n^3 terms, which reflects the fact that the naive algorithm for the product of two $n \times n$ matrices requires n^3 multiplications of scalars. It is surprising and significant that the rank of this tensor is less than n^3 . Strassen proved that when $n = 2$, its rank is at most seven, and this has led to algorithms for matrix multiplication that, for large values of n , are substantially faster than the natural one.

For further information, start with Prasolov.

In the most general sense, an algebra is a vector space V with a bilinear multiplication μ defined on it. As above we can identify μ with a cubic tensor. For v_1, \dots, v_d is a basis for V and γ_i is the element of V^* that maps a vector v to its i -th coordinate, then for x and y in V , we have

$$\mu(x, y) = \sum_{i,j=1}^d \gamma_i(x) \gamma_j(y) \mu(v_i, v_j).$$

and so we can identify μ with the element

$$\sum_{i,j} \gamma_i \otimes \gamma_j \otimes \mu(v_i, v_j)$$

of $V^* \otimes V^* \otimes V$ or, if we willing to be flexible, with an element of $V^{\otimes 3}$.

10.5 Semifields

We consider the problem of deciding which cubic tensors determine something like a field.

- (1) Suppose $T \in \text{End}(V)$ and $\text{rk}(T) = 1$. Prove that there is f in V^* and v in V such that $Tx = f(x)v$.

Chapter 11

Control

We think a linear system as a kind of ‘black box’. At time intervals $t = 1, \dots$ it receives an input, returns an output and moves to a new state. The states are elements of its *state space*, the inputs come from an *input space* and the outputs belong to the *output space*. If these elements are represented by vectors $x(i)$, $u(i)$ and $y(i)$ respectively, then they are related by the system of equations

$$\begin{aligned}x(n+1) &= Ax(n) + Bu(n), \\y(n) &= Cx(n) + Du(n),\end{aligned}$$

for all non-negative integers n . Thus the behaviour of the system is governed by the four matrices A , B , C and D , which we often write as a 2×2 matrix:

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}.$$

We will call this the *state-space* description of our system. The state-space matrix need not be square, but A must be.

What we have just described is more usually known as a discrete linear system. Since we will not consider continuous systems at any length, dropping the adjective should not cause problems.

11.1 Buffalos

By way of a first example, we consider a model for the US buffalo population, from J. J. Truxal “Introductory System Engineering”, (McGraw-Hill, New York)

1972. In this section we describe the underlying uncontrolled system; in the next section we consider the controlled version.

Let c_i and b_i respectively denote the number of female and male buffalo at the start of year i . We assume that buffalo are mature at age 2, and that each year five percent of the adults die. Female buffalo start breeding at age 2; the number of female calves born in year i is $0.12c_{i-2}$, the number of males is $0.14c_{i-2}$. Thus the population is governed by the two recurrences:

$$\begin{aligned}c_n &= 0.95c_{n-1} + 0.12c_{n-2} \\ b_n &= 0.95b_{n-1} + 0.14c_{n-2}\end{aligned}$$

We analyse the female population. If we define

$$C_n := \begin{pmatrix} c_n \\ c_{n-1} \end{pmatrix},$$

then

$$C_{n+1} = \begin{pmatrix} 0.95 & 0.12 \\ 1 & 0 \end{pmatrix} C_n.$$

Suppose

$$A := \begin{pmatrix} 0.95 & 0.12 \\ 1 & 0 \end{pmatrix}.$$

Then the minimal polynomial of A is

$$t^2 - 0.95t - 0.12,$$

which has distinct roots. Hence we can compute the spectral decomposition of A , with the result that

$$A^n = (1.0629)^n E_1 + (-0.1122)^n E_2,$$

where

$$E_1 = \begin{pmatrix} 0.9040 & 0.1021 \\ 0.8505 & 0.0960 \end{pmatrix}, \quad E_2 = \begin{pmatrix} 0.0960 & -0.1021 \\ -0.8505 & 0.9040 \end{pmatrix}.$$

(The matrices E_1 and E_2 are idempotent and $E_1 E_2 = E_2 E_1 = 0$.) From this we learn that, in the long term, the number of female buffalo will increase annually by 6.29%. The actual numbers at the end of year n will be closely approximated by the vector

$$(1.0629)^{n-1} (0.9040c_2 + 0.1021c_1).$$

This shows that even though the size of the population is sensitive to the initial conditions, even though the growth rate is not.

We now consider the males too. Suppose

$$D_n := \begin{pmatrix} c_n \\ c_{n-1} \\ b_n \end{pmatrix}.$$

Then

$$D_{n+1} = \begin{pmatrix} 0.95 & 0.12 & 0 \\ 1 & 0 & 0 \\ 0 & 0.14 & 0.95 \end{pmatrix} D_n.$$

Here the coefficient matrix is block-triangular, and its minimal polynomial is

$$(t - 0.95)(t^2 - 0.95t - 0.12).$$

- (1) Show that the male population grows as a power of 1.0629.
- (2) What is the asymptotic ratio of males to females? (It can be determined from an idempotent.)

11.2 Burgers

We continue with the model of the previous section, but we assume that each year a certain number h_n of the adult females are harvested. The equations describing the female population become

$$\begin{aligned} c_n &= 0.95c_{n-1} + 0.12c_{n-2} - h_n \\ b_n &= 0.95b_{n-1} + 0.14c_{n-2} \end{aligned}$$

which we write in matrix form as

$$C_{n+1} = AC_n - h_n \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

Let us assume that $h_n = hc_n$, for some constant h . Then we can write the resulting system as

$$C_{n+1} = A(h)C_n,$$

where

$$A(h) = \begin{pmatrix} 0.95 - h & 0.12 \\ 1 & 0 \end{pmatrix}.$$

The minimal polynomial of A_h is

$$t^2 - (0.95 - h)t - 0.12. \quad (11.2.1)$$

Given our model, we must have $0 \leq h \leq 0.95$. Let θ_h and τ_h denote the eigenvalues of $A(h)$. Then $\theta_h \tau_h = -0.12$, since this is the constant term of the minimal polynomial. It follows that θ_h and τ_h are distinct. Therefore $A(h)$ is diagonalizable, for all h . For small values of h , we see that the minimal polynomial of $A(h)$ we may assume $\theta_h \approx 1$ and τ_h is small and negative. The population will grow as a power in θ_h , and will be asymptotically constant if and only if $\theta_h = 1$. If this happens, then

$$1 - (0.95 - h) - 0.12 = 0,$$

implying that $h = 0.07$. In this case the eigenvalues are 1 and -0.12 , and idempotent corresponding to 1 is

$$\begin{pmatrix} 0.8929 & 0.1071 \\ 0.8929 & 0.1071 \end{pmatrix}.$$

- (1) Explain why the female population can stay constant when we harvest 7% of the animals annually, even though the uncontrolled growth rate is only 6.3%.

11.3 Controllability

Consider the linear system given by

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix},$$

where A is $n \times n$ and B is $n \times k$. If the initial state of the system is $x(0)$, then we have the equations

$$x(1) = Ax(0) + Bu(0)$$

$$x(2) = A^2x(0) + ABu(0) + Bu(1)$$

which leads us to the general formula

$$x(n) = A^n x(0) + \sum_{i=1}^n A^{n-i} B u(i-1).$$

Thus the state at time n is the sum of two terms, namely the state of the uncontrolled system at time n and the state of the controlled system with zero initial state. (This decomposition is an important property of linear systems.)

Define the *controllability matrix* to be

$$\mathcal{R} = (B \quad AB \quad \cdots \quad A^{n-1}B).$$

Since A is $n \times n$, its minimal polynomial has degree at most n , and so if $i \geq 0$, then A^{n+i} is a linear combination of

$$I, A, \dots, A^{n-1}.$$

Therefore the column space of \mathcal{R} is the sum of the subspaces $A^r \text{col}(B)$, where $0 \leq r < n$. It follows that if our initial state is zero, then the state of the system is always an element of $\text{col}(CM)$.

We say the pair (A, B) is *controllable* if, given any vector v in \mathbb{F}^n and starting with $x(0) = 0$, we can choose inputs $u(0), u(1), \dots, u(n-1)$ so that $x(n) = v$. We will call the system itself controllable if (A, B) is.

11.3.1 Theorem. *For a linear system, the following are equivalent:*

- (a) *The pair (A, B) is controllable.*
- (b) *The rows of the controllability matrix are linearly independent.*
- (c) *The only A -invariant subspace that contains $\text{col}(B)$ is \mathbb{R}^n .*
- (d) *No non-zero subspace of $\ker(B^T)$ is A^T -invariant.*

Proof. By the previous lemma, (a) and (b) are equivalent. The column space of the controllability matrix is the smallest A -invariant subspace that contains the columns of B , hence (b) holds if and only if (c) holds. We show that (c) and (d) are equivalent too.

Suppose $\text{rk}(\mathcal{R}) < n$. Then there is a non-zero vector f such that $f^T \mathcal{R} = 0$, and so

$$f^T A^r B = 0, \quad r = 0, 1, \dots, n-1.$$

Consequently $f^T A^r B = 0$ for all non-negative r , and therefore the A^T -invariant subspace generated by f lies in $\ker(B^T)$.

Conversely, if the A^T -invariant subspace generated by the non-zero vector f in $\ker(B^T)$ is contained in $\ker(B^T)$, that $f^T A^r B = 0$ for all r , and $\text{rk}(\mathcal{R}) < n$. \square

11.3.2 Corollary. *If B is $n \times 1$ and (A, B) is controllable, then the minimal polynomial of A has degree n .* \square

There is another concept related to controllability, sometimes called *controllability to the origin*. Suppose our system starts in some state $x(0)$ and we wish to know if there is a sequence of inputs which will drive it to the zero state.

Now the state at time r will be

$$A^r x(0) + \sum_{i=1}^r A^{r-i} B u(i).$$

Since

$$\sum_{i=1}^r A^{r-i} B u(i) \in \text{col}(\mathcal{R}),$$

we see that if there is a sequence of inputs that takes the state to zero in r steps, then $A^r x(0)$ must lie in $\text{col}(\mathcal{R})$. If $r \geq n$ and $A^r x(0) \in \text{col}(\mathcal{R})$, then there is a sequence of inputs of length r that sends the system to zero.

Thus we see, for example, that if $\text{rk } \mathcal{R} = n$, then we can bring the system to rest in n steps. To be more precise, we investigate the range of A . We note that

$$\text{col}(A^r)$$

is a nested sequence of A -invariant subspaces which is first strictly decreasing, then constant. Since $\dim(\text{col}(A)) \leq n$, it follows that when $r \geq n$,

$$\text{col}(A^r) = \text{col}(A^n).$$

We conclude that our system can be brought to rest in n steps if and only if $A^n x(0) \in \text{col}(\mathcal{R})$. Further it can be brought to rest in n steps no matter what the initial state is, if and only if

$$\text{col}(A^n) \subseteq \text{col } \mathcal{R}.$$

We conclude that any controllable system is controllable to the origin, but the latter condition is weaker.

- (1) Show that (A, B) is controllable if and only if \mathcal{R} has a right inverse.
- (2) Show that the column space of \mathcal{R} is the smallest A -invariant subspace of \mathbb{F}^n that contains the columns of B .
- (3) If A is invertible, show that (A, B) is controllable to the origin if and only if it is controllable.

11.4 Observability

Consider the linear system given by the matrix

$$\begin{pmatrix} A & B & \mathcal{R}C & D \end{pmatrix}.$$

We consider the problem of determining the initial state from the observed values of y . We have

$$x(r+1) = Ax(r) + Bu(r), \quad y(r) = Cx(r) + Du(r).$$

Since we know the values of the input vectors $u(r)$, our problem reduces to that of reconstructing $x(0)$ from the vectors $Cx(r)$. Now

$$x(r+1) = A^r x(0) + \sum_{i=1}^r A^{r-i} Bu(i);$$

since the vectors $A^{r-i} Bu(i)$ are known, the final form of our problem is to reconstruct $x(0)$ from the sequence $CA^r x(0)$ for $r = 0, 1, \dots$. Since A is $n \times n$, it follows that the first n values of this sequence determine the rest.

We say that the pair (C, A) is *observable* if the sequence

$$Cx, CAx, \dots, CA^{n-1}x$$

determines x (in all cases). The system itself is observable if (C, A) is. Define the *observability matrix* \mathcal{O} by

$$\mathcal{O} := \begin{pmatrix} C \\ CA \\ \vdots \\ CA^{n-1} \end{pmatrix}.$$

11.4.1 Theorem. *The pair (C, A) is observable if and only if the columns of the observability matrix are linearly independent.*

Proof. If the columns of \mathcal{O} are linearly independent, then it has a left inverse N . So $N\mathcal{O}x = x$, and thus we recover x . \square

11.4.2 Corollary. *The pair (C, A) is observable if and only if (A^T, C^T) is controllable.* \square

This implies for example, that (C, A) is observable if and only if no subspace of $\ker(C)$ is A -invariant.

11.5 Feedback and Controllability

Consider the system

$$x(n+1) = Ax(n) + Bu(n).$$

If we take u to be given by

$$u(n) = Kx(n) + v(n),$$

then our system becomes

$$x(n+1) = (A + BK)x(n) + Bv(n).$$

The $Kx(n)$ term is called *feedback*; the behaviour of the original system is governed by A , the behaviour of the system with feedback is governed by the matrix $A + BK$. We show that feedback does not effect controllability.

11.5.1 Lemma. *Suppose A is $n \times n$ and B is $n \times k$. Then for any $k \times n$ matrix K , the pair (A, B) is controllable if and only if $(A + BK, B)$ is.*

Proof. We show that $\text{col}(\mathcal{R}(A + BK, B)) \subseteq \text{col}(\mathcal{R}(A, B))$. Since

$$A = (A + BK) + B(-K),$$

it follows that these two column spaces are equal.

If $v \in \text{col}(\mathcal{R})$, then $Av \in \text{col}(\mathcal{R})$ and

$$BKv \in \text{col}(B) \subseteq \text{col}(\mathcal{R}),$$

whence $(A + BK)v \in \text{col}(\mathcal{R})$. It follows that $\text{col}(\mathcal{R})$ is an $(A + BK)$ -invariant subspace that contains $\text{col}(B)$, and therefore it contains the column space of $\mathcal{R}(A + BK, B)$. \square

11.5.2 Lemma. *If (A, B) is controllable and b is a non-zero column of B , then there is a matrix K such that $(A + BK, b)$ is controllable.*

Proof. Assume A is $n \times n$ and that (A, B) is controllable. We aim first to find columns b_1, \dots, b_k of B and integers r_1, \dots, r_k , such that the union of the sets

$$S(b_i, r_i) := \{b_i, Ab_i, \dots, A^{r_i-1}b_i\}$$

is a basis for \mathbb{F}^n . This is straightforward. Choose b_1 equal to b and choose r_1 to be the greatest integer such that the vectors

$$b_1, Ab_1, \dots, A^{r_1-1}b_1$$

are linearly independent. Next, assume inductively that we have found b_1, \dots, b_{j-1} and r_1, \dots, r_{j-1} such that

$$\bigcup_{i < j} S(b_i, r_i)$$

is linearly independent. The span of this set of vectors is A -invariant and so, if this set contains fewer than n vectors, there must be a column of B which it does not contain. Take b_j to be such a vector, and let r_j be the greatest integer such that the span of $S(b_j, r_j)$ contains no non-zero vectors from the span of the above union.

There is a unique linear mapping \mathcal{L} such that

$$\mathcal{L}(A^j b_i) = \begin{cases} b_{i+1}, & \text{if } j = r_i - 1; \\ 0, & \text{otherwise.} \end{cases}$$

Let L be the matrix representing \mathcal{L} . We claim that the vectors

$$b, (A+L)b, \dots, (A+L)^{n-1}b$$

are linearly independent.

Since $LA^i b_1 = 0$ if $i < r_1 - 1$ and $LA^{r_1-1} b_1 = b_2$, we see that if $i > 1$, then

$$(A+L)^{r_1-i} b = A^{r_1-i} b = A^{r_1-i} b_1$$

and

$$(A+L)^{r_1} b = A^{r_1} b_1 + b_2.$$

Starting from this, a reasonably easy induction argument, which we omit, shows that the span of the m vectors

$$(A+L)^i b, \quad i = 0, 1, \dots, m-1$$

is equal to the span of the first m vectors from

$$S(b_1, r_1) \cup \dots \cup S(b_k, r_k).$$

This proves our claim.

To complete the proof, we note that the image of \mathcal{L} is spanned by columns of B , and therefore there is a matrix K such that $L = BK$. \square

11.5.3 Corollary. *Let b be a non-zero column of B . The pair (A, B) is controllable if and only if there is a matrix K such that $(A + BK, b)$ is controllable.*

Proof. The previous lemma shows that if (A, B) is controllable and b is a non-zero column of B , then there is a matrix K such that $(A + BK, b)$ is controllable. For the converse we note that if $(A + BK, b)$ is controllable, then certainly $(A + BK, B)$ is controllable. By Lemma 11.5.1, this implies that (A, B) is controllable. \square

- (1) Let b be a non-zero element of $\text{col } B$. Show that (A, B) is controllable if and only if there is a matrix K such that $(A + BK, b)$ is controllable.

11.6 Canonical Forms

We consider first the general system

$$\begin{aligned}x(n+1) &= Ax(n) + Bu(n), \\y(n) &= Cx(n) + Du(n).\end{aligned}$$

Suppose M is invertible and $x(n) = Mz(n)$ for all non-negative n . Then we rewrite our system as

$$\begin{aligned}z(n+1) &= M^{-1}AMz(n) + M^{-1}Bu(n), \\y(n) &= MCz(n) + Du(n).\end{aligned}$$

These two systems correspond respectively to the block matrices

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}, \quad \begin{pmatrix} M^{-1}AM & M^{-1}B \\ MC & D \end{pmatrix}.$$

We say two systems related in this way are *equivalent*. We will also say that the pairs (A, B) and $(M^{-1}AM, M^{-1}B)$ are equivalent.

We now confine ourselves to the single-input case, where B is $n \times 1$. Suppose A is $n \times n$ and set $M = \mathcal{R}(A, b)$. Then

$$AM = (Ab \quad A^2b \quad \cdots \quad A^n b) = MP,$$

where F is the companion matrix of ψ_b , the minimal polynomial of A relative to b . If (A, b) is controllable, then $\text{rk}(M) = n$, and so M is invertible. It follows that

$M^{-1}AM = F$. Since $\mathcal{R}(A, b)e_1 = b$, we also find that $M^{-1}B = e_1$. We conclude that if a pair (A, b) is controllable, then our original system is equivalent to the system

$$\begin{pmatrix} F & e_1 \\ CM & D \end{pmatrix}.$$

where F is the companion matrix of the minimal polynomial of A . If we also have a single output, that is, if C is $1 \times n$, then $C = c^T$ and

$$c^T M = (c^T b \quad c^T Ab \quad \cdots \quad c^T A^{n-1} b).$$

It follows that our system is determined by the minimal polynomial of A and the entries of this vector.

From ??, we know that if F is a companion matrix of order $n \times n$, there is an symmetric invertible matrix Q such that $Q^{-1}FQ = F^T$. We see that $Qe_n = e_1$, and therefore the pair (F, e_1) is equivalent to the pair (F^T, e_n) . The pairs (C, e_1) and (C^T, e_n) are called the *controllability canonical forms* of the pair (A, b) .

There are analogous canonical forms for observable pairs (c^T, A) , but these can be deduced from our work above, applied to the controllable pair (A^T, c) .

11.7 Eigenvalues and Controllability

In this section our matrices are real matrices, but our subspaces may be complex. (For example, eigenspaces.)

11.7.1 Lemma. *The pair (A, B) is controllable if and only if the rows of $(A - \lambda I \quad B)$ are linearly independent for all complex numbers λ .*

Proof. First suppose $z \neq 0$ and

$$z^* (A - \lambda I \quad B) = 0. \tag{11.7.1}$$

Then $z^* A^r = \lambda^r z^*$ and $z^* B = 0$, so $z^* \mathcal{R} = 0$. Hence the rows of the controllability matrix are linearly dependent, and therefore (A, B) is not controllable.

On the other hand if (A, B) is not controllable, then by Theorem 11.3.1 there is an A^T -invariant subspace of $\ker(B^T)$, and this subspace must contain an eigenvector z of A^T . If the eigenvalue belonging to z is λ , then (11.7.1) is satisfied. \square

The *spectrum* of a matrix is the multiset formed by its eigenvalues, and their algebraic multiplicities. The spectrum of a real-matrix is conjugate closed—if θ is an eigenvalue, then its complex conjugate $\bar{\theta}$ is an eigenvalue with the same algebraic multiplicity.

11.7.2 Theorem. *Let A be an $n \times n$ real matrix. The pair (A, B) is controllable if and only each conjugate-closed multiset of complex numbers with size n occurs as the spectrum of some matrix $A + BK$.*

Proof. Assume first that we can choose K so that $A + BK$ has any given conjugate-closed set of complex numbers as its eigenvalues.

Suppose there is a vector z such that

$$z^T (B \quad AB \quad \cdots \quad A^{n-1}B) = 0.$$

Then, for all r and any K ,

$$z^T (A + BK)^r = z^T A^r$$

and therefore

$$z^T [(A + BK_0)^r - (A + BK_1)^r] = 0.$$

Choose K_0 so that all eigenvalues of $A + BK_0$ lie inside the unit circle, and choose K_1 so that the eigenvalues of $A + BK_1$ are the distinct n -th roots of unity. Then

$$(A + BK_1)^{ns} = I$$

for all non-negative integers s , while

$$(A + BK_0)^{ns} \rightarrow 0$$

as $s \rightarrow \infty$. It follows that $z = 0$, whence the rows of $\mathcal{R}(A, B)$ are linearly independent, and (A, B) is controllable.

We turn to the converse. We first prove the result holds in the single-input case. Suppose (A, b) is a controllable pair. We work with the equivalent canonical form (F^T, e_n) , where F^T is the transpose of the companion matrix of the minimal polynomial of A . If K is a $1 \times n$ matrix, then $e_n K$ is an $n \times n$ matrix with its first $n - 1$ rows zero, and with last row equal to K . Therefore $F^T + e_n K$ is also the transpose of a companion matrix. By varying our choice of K , we can arrange to the last row of $F^T + e_n K$ to be any desired vector, and so force

$F^T + e_n K$ to have any desired conjugate-closed set of complex numbers as its eigenvalues.

We consider the general case. Suppose (A, B) is controllable and b is a non-zero column of B . Then by Lemma 11.5.2 there is a matrix K such that $(A + BK, b)$ is controllable. By what we have just proved, for each conjugate-closed set of complex numbers, there is a $1 \times n$ matrix K_1 such that

$$A + BK + bK_1$$

has this set as its eigenvalues. But $b = Be_r$ for some r , and so

$$BK + bK_1 = BK + Be_r K_1 = B(K + e_r K_1),$$

and so our result is proved. □

11.8 Observers

Consider the discrete dynamical system given by the equations

$$\begin{aligned} x(n+1) &= Ax(n) + Bu(n) \\ y(n) &= Cx(n) + Du(n). \end{aligned}$$

We want to construct a second system which will accept both the input and the output of the first system as its inputs, and as produce as its own output at least an approximation to the state of our first system. To construct such a system, we consider a second system based on the one above:

$$\begin{aligned} \hat{x}(n+1) &= A\hat{x}(n) + Bu(n) + L(y(n) - \hat{y}(n)) \\ \hat{y}(n) &= C\hat{x}(n) + Du(n). \end{aligned}$$

If this system has the property that $x(n) - \hat{x}(n) \rightarrow 0$ as $n \rightarrow \infty$, we call it an *asymptotic observer*. If A is $m \times m$ and $x(n) = \hat{x}(n)$ when $n > m$, we call it an *exact observer*. The choice of L is up to us. We calculate

$$\begin{aligned} x(n+1) - \hat{x}(n+1) &= A(x(n) - \hat{x}(n)) - L(y(n) - \hat{y}(n)) \\ &= A(x(n) - \hat{x}(n)) - L(Cx(n) - C\hat{x}(n)) \\ &= (A - LC)(x(n) - \hat{x}(n)). \end{aligned}$$

Now the initial difference $x(0) - \hat{x}(0)$ can be any vector, so we conclude that we have an asymptotic observer if and only if $A - LC$ is a contraction, and an exact observer if and only if $A - LC$ is nilpotent.

If A is $n \times n$ and B is $n \times m$, we say that the pair (A, B) is *stabilizable* if there is an $m \times n$ matrix K such that all eigenvalues of $A + BK$ lie inside the unit circle. (In other terms, $A + BK$ is a contraction.) Every controllable pair is stabilizable. If C is $\ell \times n$, we say that (C, A) is *detectable* if (A^T, C^T) is stabilizable or, equivalently, if there is a matrix L such that $A + LC$ is a contraction.

11.8.1 Theorem. *An asymptotic observer exists if and only if (C, A) is detectable. An observer exists if and only if (C, A) is observable. \square*

11.9 Transfer Matrices

We introduce a very important tool in the study of discrete dynamical systems: transfer matrices.

We first present this in a special case, coming from coding theory. We suppose that a sequence $(u_i)_{i \geq 0}$ of binary vectors is encoded by a device as a second sequence $(y_i)_{i \geq 0}$ of binary vectors. In the simplest case, we have a matrix D and u_i is mapped to Du_i . But we are going to assume that our device has a state x_i (another binary vector) and that y_i is computed according to the system

$$x_{i+1} = Ax_i + Bu_i \tag{11.9.1}$$

$$y_i = Cx_i + Du_i. \tag{11.9.2}$$

Here A, B, C and D are binary matrices and A is square. (For a coding theorist it might be natural to assume D is $n \times k$; the matrix A is square.) The first problem that arises is to reconstruct the inputs u_i given the outputs y_i (and the four matrices A, B, C, D). In the real applications, the vectors y_i are corrupted by noise, and we also have the harder task of first determining the uncorrupted values of the outputs.

We say that the system described by the four matrices is a *convolutional encoder*. The space of possible output sequences is a *convolutional code*. Convolutional codes are important in practice.

To make further progress, we introduce generating functions. A convolutional encoder takes an input sequence

$$u_0, u_1, u_2, \dots$$

and converts it to an output sequence

$$y_0, y_1, y_2, \dots$$

In any practical situation, the vectors u_i will be zero for all sufficiently large i , but we defer imposing this as a requirement. One standard way to deal with infinite sequences is to encode them as formal power series, and so we define

$$U(z) := \sum_{i \geq 0} z^{-i} u_i, \quad Y(z) := \sum_{i \geq 0} z^{-i} y_i.$$

These can be viewed as formal power series in the variable z^{-1} with vectors as coefficients, or as vectors whose entries are formal power series of \mathbb{F} . (We tend to prefer the latter view.) We say that $U(z)$ is a generating function for the sequence $(u_i)_{i \geq 0}$.

Next we assume that $x_0 = 0$ and introduce the generating function $X(z)$. The defining equations for our encoder give us

$$zX(z) = AX(z) + BU(z), \quad Y(z) = CX(z) + DU(z),$$

and consequently

$$Y(z) = (D + C(zI - A)^{-1}B)U(z).$$

It follows that our encoder is completely specified by the proper rational matrix

$$G(z) := D + C(zI - A)^{-1}B.$$

If we have a discrete dynamical system over a field \mathbb{F} , given by the matrix

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix} \tag{11.9.3}$$

we define the *transfer matrix* of the system to be

$$D + C(zI - A)^{-1}B.$$

The transfer matrix completely determines the response of our system to a given input sequence (given that $x_0 = 0$). If e_i denotes the i -th standard basis vector, then the generating function of the output sequence corresponding to the input sequence

$$e_i, 0, 0, \dots$$

is the i -th column of $G(z)$. This provides a very natural interpretation of the columns of $G(z)$, and shows that we can find the transfer matrix of a system by determining its response to each of the above input sequences. In particular it is not unusual to be given the transfer matrix of a system, rather than the state-space description.

It may seem more natural to use formal power series in z rather than z^{-1} , but the above choice is standard in control theory.

Chapter 12

The Smith Normal Form

In this chapter we study some linear algebra over rings. The most important rings we use are \mathbb{Z} and $\mathbb{F}[x]$.

12.1 Domains

Let R be a commutative ring. We say that an element a of R divides an element b if $b = ax$ for some x . We call R a *domain* if it has no divisors of zero, that is, if $a, b \in R$ and $ab = 0$ then $a = 0$ or $b = 0$. Clearly any field is a domain. Further examples are provided by the integers \mathbb{Z} and $\mathbb{F}[x]$, the ring of polynomials in x with coefficients from \mathbb{F} .

An ideal of R is a non-empty subset I such that if $a \in I$ and $r \in R$, then $ra \in I$. The even integers form an ideal in \mathbb{Z} . The polynomials p in $\mathbb{F}[x]$ such that $p(1) = 0$ provide a second example. If I and J are subsets of R , then IJ is given by

$$IJ := \{ab : a \in I, b \in J\}.$$

Thus the subset I of R is an ideal if $RI \subseteq I$. The only ideal of R that contains 1 is R itself. It follows that a proper ideal cannot contain an invertible element of R . If $S \subseteq R$, then the set SR is an ideal; we call it the ideal generated by S . It consists of all R -linear combinations of the elements of S . An ideal generated by a single element is called a *principal ideal*. For example, the even integers $2\mathbb{Z}$ form a principal ideal in \mathbb{Z} . If I is the principal ideal generated by d , then I consists of the elements of R that are divisible by d . A *principal ideal domain* is a ring in which every ideal is principal. Both \mathbb{Z} and $\mathbb{F}[x]$ are examples.

An ideal I is *prime* if it is a proper ideal and, whenever $ab \in I$, either a or b lies in I . If $m \in \mathbb{Z}$, then $m\mathbb{Z}$ is a prime ideal if and only if m is a prime. A ring is a domain if and only if the zero ideal is prime.

Suppose R is a principal ideal domain and $a, b \in R$. The ideal generated by a and b is generated by some element d , which divides both a and b . Since this ideal consists of the R -linear combinations of a and b , there are elements r and s of R such that

$$d = ra + sb.$$

It follows that if c divides a and b , then c divides d and therefore d is a greatest common divisor of a and b .

If d divides e and e divides d , we have

$$d = d_1 e, \quad e = e_1 d$$

whence $d = d_1 e_1 d$. Therefore

$$(1 - d_1 e_1) d = 0$$

and so $d_1 e_1 = 1$; hence both d_1 and e_1 are units of R . It follows that, in a principal ideal domain, any two non-zero elements have a greatest common divisor, which is unique up to multiplication by a unit.

It can be difficult to verify that a given ring is a principal ideal domain. There is one case where it is easy. We say R is a *Euclidean domain* if there is a function ρ from $R \setminus 0$ to \mathbb{N} such that

- (a) If $a, b \in R$ then $\rho(ab) \geq \rho(a)$.
- (b) If $a, b \in R$, there are elements q and r such that $b = qa + r$ and $\rho(r) < \rho(a)$.

The advantage of Euclidean domains is that we can compute the greatest common divisor of any two elements using the usual Euclidean algorithm. Also, a Euclidean domain is a principal ideal domain.

We consider examples. If $R = \mathbb{Z}$, take $\rho(a)$ to be $|x|$. If $R = \mathbb{F}[x]$, use $\rho(p) = \deg(p)$. If $p, q \in \mathbb{F}[x]$, we say the rational function p/q is *proper* if $\deg(p) < \deg(q)$. If we define ρ by

$$\rho\left(\frac{p}{q}\right) := \deg(q) - \deg(p).$$

The set of proper rational functions over \mathbb{F} , with this function ρ , forms a Euclidean domain. If $\mathbb{F} = \mathbb{C}$, the strictly proper rational functions are the rational functions that are bounded at infinity.

1. Prove that a finite domain is a field.
2. If (R, ρ) is a Euclidean domain and x is a unit in R , show that $\rho(ax) = \rho(a)$ for all a in R .

12.2 Localization

Let R be a domain. A subset S of R is *multiplicatively closed* if

- (a) $0 \notin S$ and $1 \in S$,
- (b) If $a, b \in S$, then $ab \in S$.

In \mathbb{Z} , the set of integers not divisible by a given prime is multiplicatively closed. The set of non-zero elements of R is also multiplicatively closed.

Using S , we can construct a new ring, denoted $R[S^{-1}]$. Its elements are equivalence classes of ordered pairs from $R \times S$. We define (a, s) and (b, t) to be equivalent if there is an element x of R such that $b = ax$ and $t = sx$. The product of the pairs (a, s) and (b, t) is (ab, st) ; their sum is

$$(at + bs, st).$$

These definitions will seem more familiar if we write our pairs as ratios a/s . We then see that if $R = \mathbb{Z}$ and S consists of the non-zero integers, $R[S^{-1}] = \mathbb{Q}$. If $R = \mathbb{F}[x]$ and S consists of all powers of x , then $R[S^{-1}]$ is known as the ring of *Laurent polynomials*. It consists of the rational functions of the form $x^k p(x)$, where $p \in \mathbb{F}[x]$ and $k \in \mathbb{Z}$.

If $S = R \setminus \{0\}$, then the ring $R[S^{-1}]$ is called the *quotient field* of R . The quotient field of \mathbb{Z} is \mathbb{Q} , as we have just noted. The quotient field of $\mathbb{F}[x]$ is ring of rational functions in x , denoted $\mathbb{F}(x)$.

We can view $R[S^{-1}]$ as being constructed by adjoining the multiplicative inverse of each element of S to R . The elements of $R[S^{-1}]$ of the form $a/1$ form a subring isomorphic to R . If $a \in R$ and $s \in S$, then $a/1$ and a/s generate the same ideal. It follows from this that ideals of $R[S^{-1}]$ correspond to the ideals I of R such that $I \cap S = \emptyset$. An important consequence is that $R[S^{-1}]$ is a principal ideal domain if R is.

We consider some examples. Let $R = \mathbb{C}[x]$ and let C be a subset of \mathbb{C} , for example, the unit disc. Then the polynomials $p(x)$ with no zeros in C form a multiplicatively closed subset S . The ring $R[S^{-1}]$ consists of the rational functions with no pole in the unit disc.

1. Prove that if I is an ideal of R , then $R \setminus I$ is multiplicatively closed if and only if I is prime.
2. Let S be a multiplicatively closed subset of the domain R . Prove that each ideal of $R[S^{-1}]$ consists of the elements a/s , where a comes from a given ideal I of R , and $s \in S$.

12.3 Fitting

Fitting Invariants

Let A be an $m \times n$ matrix over a ring R . We define the *Fitting invariant* $F_k(A)$ to be the ideal generated by the $k \times k$ minors of A , where $1 \leq k \leq \min\{m, n\}$. Thus

$$F_k(A) = F_1(\Lambda_k(A)).$$

If R is a principal ideal domain, then the ideal $F_k(A)$ is generated by an element f_k , and so we may use the sequence of elements $f_1, \dots, f_{m \wedge n}$, rather than the ideals $F_k(A)$.

12.3.1 Lemma. *Let A be an $m \times n$ matrix over R , where $m \leq n$. Then the following are equivalent:*

- (a) A has a right inverse.
- (b) $F_k(A) = R$ for $k = 1, \dots, m$.
- (c) $F_m(A) = R$.

Proof. First suppose that B is a right inverse for A over R . Then B is $m \times n$ and, since $AB = I_m$, we have

$$\Lambda_k(A)\Lambda_k(B) = I_{\binom{m}{k}}.$$

From this we see that 1 is an R -linear combination of elements of each row of $\Lambda_k(A)$, and consequently $F_k(A) = R$.

Now assume that $F_m(A) = R$. Let S be a set of m columns of A and set d_S equal to $\det A_S$. Let M be the $n \times m$ matrix such that $M_S = \text{adj}(A_S)$ and $Me_i = 0$ if $i \notin S$. Then

$$AM = d_S I.$$

If T is a second subset of m columns of A and N is constructed analogously to M , then $AN = d_T I$ and therefore

$$A(xM + yN) = (xd_S + yd_T)I.$$

It follows that if the minors d_S generate R , then there is a right inverse for A . \square

One consequence of this lemma is that if A is $m \times n$ and $F_m(A) = R$, then $F_k(A) = R$ for $k = 1, \dots, m$. We recall that two matrices A and B are equivalent over R if there are invertible matrices P and Q over R such that $B = PAQ$.

12.3.2 Lemma. *Let A and B be two $m \times n$ matrices over R . If A and B are similar, they have the same Fitting invariants.*

Proof. By Lemma 12.3.1 it follows that $F_k(P) = R$ for all k . Now

$$\Lambda_k(PA) = \Lambda_k(P)\Lambda_k(A)$$

and it follows that $F_k(PA) \subseteq F_k(A)$. Applying the same argument to the pair of matrices P^{-1} and PA , we get

$$F_k(A) = F_k(P^{-1}(PA)) \subseteq F_k(PA).$$

Accordingly $F_k(PA) = F_k(A)$, as claimed.

In Section 12.5, we will see that if R is a principal ideal domain, then two matrices of the same order are equivalent if and only if they have the same Fitting invariants. Note that A and A^T have the same Fitting invariants.

1. Let A be a $m \times n$ matrix over \mathbb{Z} . Show that if, for each prime p , the rank of A modulo p equals its rank over \mathbb{Q} , then the greatest common divisor of the $m \times m$ minors of A is 1.

12.4 Hermite

Suppose A is a matrix over a ring R . There are three *elementary row operations* we can apply to A :

- (1) Multiply a row by a unit.
- (2) Swap two rows.

(3) Add a multiple of one row to another.

Similarly we have three types of elementary column operations. An elementary matrix is a matrix we get by applying an elementary row operation to the identity. (It makes no difference if we use elementary column operations.) An elementary matrix is invertible, and its inverse is elementary. Over a field, every invertible matrix is a product of elementary matrices. Over rings the situation is much more complicated, but we can deal with Euclidean domains.

The main tool is an analog of reduced row echelon form. The key observation is the following. Suppose $a, b \in R$ and that d is a greatest common divisor of a and b . Then there are elements s and t in R such that $sa + tb = d$. Suppose $a = a_1d$ and $b = b_1d$. Then

$$\begin{pmatrix} s & t \\ -b_1 & a_1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} d \\ 0 \end{pmatrix}.$$

Note the determinant of the 2×2 matrix on the left is 1, and so it is invertible. We can use this to show that if A is a matrix over R , then there is an invertible matrix F such that FA is in row echelon form.

If we are working over a Euclidean domain, there are two refinements. Suppose R is Euclidean relative to the function ρ . Then by applying elementary row operations to the row echelon form of A , we may convert it to a matrix B in row echelon form such that, if the first non-zero element of row j is in column k , then

$$\rho(B_{i,k}) < \rho(B_{j,k}).$$

We say that B is in *Hermite normal form*.

The second refinement is the following result.

12.4.1 Lemma. *Suppose that R is a Euclidean domain. If a and b are elements of R with a greatest common divisor d , there is a product E of elementary matrices such that*

$$E \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} d \\ 0 \end{pmatrix}.$$

Proof. We prove this by using the Euclidean algorithm. First, premultiplying by a permutation matrix if needed, we may assume that $\rho(a) \geq \rho(b)$. Then

$$a = q_1b + r_1,$$

where $\rho(r_1) < \rho(b)$. We have

$$\begin{pmatrix} 0 & 1 \\ 1 & -q_1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} b \\ r_1 \end{pmatrix}.$$

Now

$$\begin{pmatrix} 0 & 1 \\ 1 & -q_1 \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & -q_1 \\ 0 & 1 \end{pmatrix}$$

and so we have converted $\begin{pmatrix} a \\ b \end{pmatrix}$ to $\begin{pmatrix} b \\ r_1 \end{pmatrix}$ by premultiplying it by two elementary matrices. If $r_1 = 0$, we are done. If $r_1 \neq 0$, there is an element q_2 in R such that

$$\begin{pmatrix} 0 & 1 \\ 1 & -q_2 \end{pmatrix} \begin{pmatrix} b \\ r_1 \end{pmatrix} = \begin{pmatrix} r_1 \\ r_2 \end{pmatrix}$$

and $\rho(r_2) < \rho(r_1)$. Thus each step of the Euclidean algorithm is equivalent to multiplying an element of R^2 by two elementary matrices, and the lemma follows. \square

We summarise our conclusions.

12.4.2 Theorem. *If A is a matrix over a Euclidean domain R , it can be converted to Hermite normal form by elementary column operations.* \square

1. Show that over a Euclidean domain, each invertible matrix is a product of elementary matrices.
2. Suppose B and C are $m \times n$ matrices over a Euclidean domain in Hermite normal form, with linearly independent columns. If $B = CG$, show that G is diagonal. (Hint: first show that G is lower triangular.)

12.5 Smith Normal Form

Let A and B be two $m \times n$ matrices over a commutative ring R . (Think $\mathbb{Z}[x]$.) We say that A and B are *equivalent* over R if there are invertible matrices P and Q such that $PAQ = B$. We want to decide if two given matrices are equivalent.

12.5.1 Theorem. *Let A be a matrix over a principal ideal domain R . Then there is a unique matrix D over R which is equivalent to A such that $D_{ij} = 0$ if $i \neq j$ and $D_{i,i}$ divides $D_{i+1,i+1}$ for $i = 1, \dots, n-1$.*

Proof. Although it is not strictly necessary, we first show that A is equivalent to a matrix D such that $D_{ij} = 0$ if $i \neq j$, and only then show that D can be arranged to have the form stated.

Suppose a and b are two elements of R , and suppose that the ideal they generate is generated by d . Then there must be elements s and t of R such that $sa + tb = d$. Further, there are elements a_1 and b_1 such that $a = a_1 d$ and $b = b_1 d$. Hence

$$\begin{pmatrix} s & t \\ -b_1 & a_1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} d \\ 0 \end{pmatrix}.$$

As $sa_1 + tb_1 = 1$ the determinant of

$$\begin{pmatrix} s & t \\ -b_1 & a_1 \end{pmatrix}$$

is 1 and therefore this matrix is invertible.

If the i -th of A is x and the j -th row is y and we may replace x by $sx + ty$ and y by $-b_1 x + a_1 y$, the resulting matrix is equivalent to A .

We may permute the columns of A so that any zero columns are last. Having done this, we may convert A to an equivalent matrix where $A_{1,1} = a \neq 0$ and $A_{i,1} = 0$ if $i > 1$. If a divides each entry of the first row of A then A is equivalent to a matrix of the form

$$\begin{pmatrix} a & 0 \\ 0 & A_1 \end{pmatrix}$$

and we can prove our claim by induction.

If a does not divide each entry in the first row, then we may operate on the columns of A , converting it to an equivalent matrix with $A_{1,1} = a'$ and $A_{1,j} = 0$ if $j > 1$. Further the ideal generated by a is properly contained in the ideal generated by a' . We hope now that a' divides each entry in the first column of A . If so then we reduce to the previous induction. If not, we operate on the rows again. Since at each stage the ideal generated by $A_{1,1}$, and since R does not contain an infinite increasing sequence of ideals, we conclude that A is equivalent to a matrix with $A_{i,1} = 0$ when $i > 1$ and $A_{1,j} = 0$ when $j > 1$. This proves our claim.

To reduce R to the required form, we observe that the two matrices

$$\begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix}, \quad \begin{pmatrix} a & 0 \\ sa + tb & b \end{pmatrix}$$

are equivalent; given this it is easy to see R is equivalent to a matrix satisfying the divisibility condition we gave.

The problem left is to prove that R is unique. This follows because A and D have the same Fitting ideals, and because a diagonal matrix which satisfies our divisibility condition is determined by its Fitting ideals. \square

If R is a Euclidean domain, then we can use elementary row and column operations rather than the 2×2 matrices we described.

The matrix R whose existence is guaranteed by the theorem is called the *Smith normal form* of A . If A is square then $\det(A)$ is a unit times $\det(R)$. Computing the Smith normal form, even over \mathbb{Z} , is one of the more difficult problems in linear algebra. If implemented as described then the number of digits in an entry can double at each step.

Generally one only meets the Smith normal form for matrices over \mathbb{Z} and over $\mathbb{F}[z]$; there are a number of interesting Euclidean domains that arise in control theory, related to rational functions. Call a rational function p/q in $\mathbb{F}(z)$ *bounded* if $\deg p \leq \deg q$. If we define

$$\rho_1(p/q) = \deg q - \deg p,$$

the bounded rational functions form a Euclidean domain relative to the function ρ_1 .

For a second example, let S be a subset of the complex plane, and call a polynomial *stable* if its zeros all lie in S . The set of stable polynomials is multiplicatively closed, and so the rational functions p/q where q is stable form a ring. If we define $\rho_2(p/q)$ to be the number of zeros of p not in S then this ring is a Euclidean domain relative to ρ_2 .

The intersection of these two rings has the baroque denotation RH_+^∞ . If, as is standard, S is the open left half-plane, this ring consists of the rational functions that are uniformly bounded on the closed right half-plane. It is a Euclidean domain relative to the function $\rho_1 + \rho_2$.

Chapter 13

Polynomial and Rational Matrices

A *polynomial matrix* is a matrix whose entries come from the ring $\mathbb{F}[z]$. A *rational matrix* is a matrix whose entries come from the field of rational functions $\mathbb{F}(z)$. We will also have occasion to consider matrices whose entries are formal power series or Laurent series, but we will not assign names to these. Any matrix polynomial $A(z)$ can be written as a polynomial in z with coefficients A_i from $\text{Mat}_{m \times n}(\mathbb{F})$:

$$A(z) = \sum_i A_i z^i.$$

(This encodes an isomorphism between the ring of polynomial matrices, and the ring of polynomials with matrix coefficients, which goes beyond the level of sophistication to which we aspire.) The *degree* of a matrix polynomial is the maximum degree of an entry. We will also be concerned with the degrees of the rows and/or columns of polynomial matrices. The key here is to note that each column of a polynomial matrix is a polynomial matrix, and so has a well-defined degree.

We consider one pertinent example. If A is $n \times n$, then $zI - A$ is a polynomial matrix with degree one. We have

$$(zI - A) \text{adj}(zI - A) = \det(zI - A)I.$$

Here $\text{adj}(zI - A)$ is also a matrix polynomial, with degree $n - 1$, and

$$(zI - A)^{-1} = \frac{1}{\det(zI - A)} \text{adj}(zI - A).$$

In this chapter, we will study the basic properties of polynomial and rational matrices.

13.1 Series

A rational function is *proper* if its numerator and denominator have the same degree, and *strictly proper* if the degree of its numerator is less than the degree of its denominator. A rational matrix is proper if its elements are proper and strictly proper if they are strictly proper. The proper rational matrices form a ring, and the strictly proper rational matrices form a proper ideal in this ring.

We can view the ring of polynomials $\mathbb{F}[z]$ as a subring of the ring of formal power series $\mathbb{F}[[z]]$. This has some use, for example if $p(z)$ is a polynomial and $p(0) \neq 0$, then $p(z)$ has a multiplicative inverse in $\mathbb{F}[[z]]$. In a similar way, we can represent rational functions by formal Laurent series.

Suppose

$$p(z) = z^n + p_1 z^{n-1} + \cdots + p_n.$$

Then

$$p(z)^{-1} = z^{-n} \left(1 + \frac{p_1}{z} + \cdots + \frac{p_n}{z^n} \right)^{-1}$$

Hence $p(z)^{-1}$ has a formal power series expansion in z^{-1} , and it follows that any rational function has an expansion as a formal Laurent series in z^{-1} . If $p(z)/q(z)$ is a rational function then

$$\frac{p(z)}{q(z)} = \sum_{i=-k}^{\infty} a_i z^{-i},$$

where $k = \deg(p) - \deg(q)$. Hence the ring of rational functions in z is isomorphic to a subring of the ring of Laurent series in z^{-1} , and the image of the proper rational functions under this isomorphism is the ring of formal power series in z . The strictly proper rational functions map to the formal power series with constant term equal to 0.

Since we have used nothing more than the geometric series expansion, everything goes over to matrix rational functions: these are isomorphic to a subring of the ring of Laurent series in z^{-1} with matrix coefficients, proper rational matrices correspond to formal power series and strictly proper rational matrices to formal power series with constant term equal to 0. From this we see, for example, that the proper rational matrices form a ring, and the strictly proper rational functions form an ideal in this ring. We note one other property we will need.

13.1.1 Lemma. *If $M(z)$ is a strictly proper rational matrix, then $I + M(z)$ is invertible, and its inverse is a proper rational matrix.*

Proof. Since $M(z)$ is strictly proper it has a series expansion

$$M(z) = \sum_{i \geq 0} M_i z^{-i}$$

Hence $I + M(z)$ is a formal power series with constant term I , and therefore it has a multiplicative inverse, which is again a formal power series with constant term I . \square

13.2 Polynomial Matrices

We develop some of the basic properties of polynomial matrices.

Every polynomial matrix is a rational matrix. Since

$$A(z) \operatorname{adj}(A(z)) = \det(A(z))I$$

we see that if $\det(A(z)) \neq 0$, then

$$\frac{1}{\det(A(z))} \operatorname{adj}(A(z))$$

is the inverse of $A(z)$ in the ring of rational matrices. Thus a polynomial matrix $A(z)$ has a rational inverse if and only if its determinant is not zero, although $A(z)$ may not be invertible for certain values of z in \mathbb{F} . A polynomial matrix has a polynomial inverse if and only if its determinant is a non-zero constant. We say that a square matrix over a ring is *unimodular* if its determinant is a unit. Since the units in $\mathbb{F}[x]$ are the non-zero constants, a polynomial matrix is unimodular if and only if it has a polynomial inverse. More generally, we recall from Section 12.3 that an $m \times n$ matrix A over a ring R has a right inverse if and only if the ideal generated by the $m \times m$ minors of A is equal to R .

Suppose $A(z)$ is an $m \times n$ polynomial matrix with linearly independent rows that is not right invertible. Then the greatest common divisor of the $m \times m$ minors of $A(z)$ is a polynomial of positive degree. It follows that there are values of z in the algebraic closure of \mathbb{F} such that $\operatorname{rk}(A(z)) < m$.

We may write any $m \times n$ polynomial matrix $A(z)$ in the form

$$A(z) = HS(z) + L(z),$$

where $S(z)$ is diagonal with i -th diagonal entry equal to z^{d_i} and $L(z)$ is a matrix whose column degrees are each less than the column degrees of $A(z)$. We call H

the *leading coefficient matrix* of $A(z)$, and we define $A(z)$ to be *column reduced* if the columns of H are linearly independent.

By way of example, if A is square matrix over \mathbb{F} , then $zI - A$ is column reduced. If $A(z)$ is invertible and in Hermite normal form, then $H = I$ and so $A(z)$ is column reduced. A square polynomial matrix

$$\sum_{i=0}^d A_i z^i$$

is said to be *regular* if A_d is invertible. A regular polynomial matrix is column reduced, but the converse fails. (Find an example.)

13.2.1 Lemma. *Let $A(z)$ be a polynomial matrix with linearly independent columns. Then there is a product of elementary matrices $E(z)$ such that $A(z)E(z)$ is column reduced, and the degree of each column of $A(z)E(z)$ is no greater than that of the corresponding column of $A(z)$.*

Proof. Suppose some set of columns of H is linearly dependent. Choose a minimal such subset C , and from this choose a column, i say, with largest possible degree. Then the i -th column is a linear combination of the remaining columns in C , and so there is a product $E(z)$ of elementary matrices with determinant 1 such that the degree of the i -th column of $A(z)E(z)$ is less than the degree of the i -th column of $A(z)$, and all other columns have the same degree in both matrices. We may continue reducing the degrees of columns in this way, until we reach a matrix whose leading coefficient matrix has full rank. \square

13.2.2 Lemma. *Let $A(z)$ be an $n \times k$ polynomial matrix, where $n \geq k$, and let d_i be the degree of its i -th column. Then the degree of a $k \times k$ minor of $A(z)$ is at most $\sum_i d_i$ and equality holds if and only if $A(z)$ is column reduced. \square*

We follow the common convention that the degree of the zero polynomial is $-\infty$.

13.2.3 Lemma. *Let $A(z)$ be a column-reduced $m \times n$ polynomial matrix and let d_i be the degree of its i -th column. Let $p(z)$ be a vector of polynomials of length n , with i -th component $p_i(z)$. Then*

$$\deg(A(z)p(z)) = \max_i \{d_i + \deg(p_i(z))\}.$$

Proof. We write

$$A(z) = HS(z) + L(z),$$

where H is the leading coefficient matrix of $A(z)$ and $S(z)$ is diagonal with i -th entry z^{d_i} . Then

$$A(z)p(z) = HS(z)p(z) + L(z)p(z).$$

The i -th entry of $S(z)p(z)$ is $z^{d_i}p_i(z)$, and consequently

$$\deg HS(z)p(z) \leq \deg(S(z)p(z)) = \max_i \{d_i + \deg(p_i(z))\}. \quad (13.2.1)$$

The degree of the i -th entry of $L(z)p(z)$ is less than $d_i + \deg(p_i(z))$, which proves that the right side of (13.2.1) is an upper bound on the degree of $A(z)p(z)$. Since the columns of H are linearly independent, it has a left inverse, K say. Then

$$KHS(z)p(z) = S(z)p(z),$$

from which we see that

$$\deg(HS(z)p(z)) \geq \deg(S(z)p(z)).$$

Hence $\deg(HS(z)p(z)) = \deg(S(z)p(z))$, and the theorem follows. \square

13.2.4 Theorem. *Suppose $A(z)$ and $B(z)$ are column-reduced polynomial matrices, with columns arranged in increasing order of degree. If $C(z)$ is unimodular and $A(z) = B(z)C(z)$, then $A(z)$ and $B(z)$ have the same column degrees and $C(z)$ is upper triangular.*

Proof. Assume $A(z)$ and $B(z)$ are $m \times n$ and suppose that $A(z)$ has column degrees d_1, \dots, d_n and $B(z)$ has column degrees e_1, \dots, e_n . If $p(z)$ is the r -th column of $C(z)$, then the degree of $B(z)p(z)$ is at least the maximum of the degrees d_i such that $p_i(z) \neq 0$. It follows that $d_1 \geq e_1$. Since $C(z)$ is unimodular, $C(z)^{-1}$ is polynomial and since $B(z) = A(z)C(z)^{-1}$, we also see that $e_1 \geq d_1$. Hence $d_1 = e_1$ and $C(z)_{i,1} = 0$ if $i > 1$.

Now let $A_1(z)$ and $B_1(z)$ be the matrices we get by deleting the first column from $A(z)$ and $B(z)$ respectively, and let $C_1(z)$ be the matrix we get by deleting the first row and column from $C(z)$. Then A_1 and B_1 are reduced and C_1 is unimodular, and the theorem follows by induction on m . \square

(1) Prove that a column-reduced matrix has a rational left inverse.

13.3 Paraunitary Matrices

Suppose

$$A(z) = \sum_{r=0}^m A_r z^r,$$

where $A_r \in \text{Mat}_{n \times n}(\mathbb{C})$ and let $A^*(z^{-1})$ be given by

$$A^*(z) = \sum_{r=0}^m A_r^* z^{-r}.$$

We say that $A(z)$ is paraunitary if

$$A(z)A^*(z^{-1}) = I.$$

One consequence of this definition is that if $A(z)$ is paraunitary, then $A(z)$ is unitary when $\|z\| = 1$. The product of paraunitary matrices is paraunitary.

By way of example if $v \in \mathbb{C}^n$ and $\|v\| = 1$, then an easy computation shows that

$$V(z) := I - vv^* + zvv^*$$

is paraunitary. We also see that $V(1) = I$ and, with some effort, that

$$\det V(z) = 1.$$

Paraunitary matrices of this type are called *primitive*. Note that $V^*(z) = V(z)$, so $V(z)V(z^{-1}) = I$.

13.3.1 Lemma. 3.1 *If $A(z)$ is paraunitary, then $\det A(z) = z^m$ for some non-negative integer m .*

Proof. Suppose $p(z) := \det A(z)$, and let $\bar{p}(z)$ be the polynomial whose coefficients are the complex conjugates of those of $p(z)$. Then $\bar{p}(z^{-1}) = \det A^*(z^{-1})$ and since $A(z)A^*(z^{-1}) = I$, we have

$$p(z)\bar{p}(z^{-1}) = \det(A(z))\det(A^*(z^{-1})) = 1.$$

Suppose $p(z)$ has degree d and that z^e is the highest power of z that divides $p(z)$. Then

$$p(z)\bar{p}(z^{-1}) = \frac{p(z)q(z)}{z^d},$$

where $q(z)$ is a polynomial of degree $d - e$. Therefore $p(z)q(z)$ has degree $d - e$, and the lemma follows at once. \square

13.3.2 Lemma. *Let $A(z)$ be a paraunitary matrix. Then $A(z)$ is constant if and only if $\det A(z) = 1$.*

Proof. Suppose

$$A(z) = \sum_{r=0}^m z^r A_r$$

and $A_m \neq 0$. If $m = 0$ then $A(z)$ is constant and $\det A(z) = 1$.

If $m > 0$, then the coefficient of z^{-m} in the product $A(z)A^*(z^{-1})$ is $A_0A_m^*$, whence $A_0A_m^* = 0$ and A_0 is singular. Then

$$\det A(z) \det(A_0 + zB(z)),$$

where $B(z)$ is a polynomial matrix. Therefore the constant term of $\det A(z)$ is $\det A(0)$, which is zero. We conclude that $\det A(z)$ is a positive power of z . \square

13.3.3 Theorem. *If $A(z)$ is a paraunitary matrix and $\det A(z) = z^d$, then $A(z) = A(1)W(z)$, where $W(z)$ is the product of d primitive paraunitary matrices.*

Proof. If $d = 0$ then $A(z) = A(1)$ and there is nothing to prove, so we assume $d > 0$. As in the proof of the previous lemma, it follows that the constant term A_0 in $A(z)$ is singular and therefore there is a unit vector v such that $v^*A_0 = 0$.

Suppose $A(z) = \sum_{r=0}^m A_r z^r$ and

$$V(z) := I - vv^* + zvv^*$$

and consider the product

$$B(z) = V(z^{-1})A(z) = (I - vv^* + z^{-1}vv^*)(A_0 + A_1z + \cdots + A_mz^m).$$

Since $v^*A_0 = 0$, we see that $B(z)$ is a polynomial matrix and hence that is paraunitary.

Since $V(z)B(z) = A(z)$ we have

$$z \det B(z) = z^d$$

and consequently $\det B(z) = z^{d-1}$. The theorem follows now by induction on the degree of $\det A(z)$. \square

Paraunitary matrices play a significant role in the theory of filter banks and in some treatments of wavelets. For the latter, see Resnikoff and Wells “Wavelet Analysis” (Springer, New York) 1998.

13.4 Division

If $a(z)$ and $d(z)$ are polynomials over a field, there are unique polynomials $q(z)$ and $r(z)$ such that $\deg r < \deg d$ and

$$a(z) = q(z)d(z) + r(z).$$

We establish a matrix version of this.

13.4.1 Theorem. *Suppose $D(z)$ and $N(z)$ are polynomial matrices of orders $n \times n$ and $m \times n$ respectively, and that $D(z)$ is column reduced. Then $N(z)D(z)^{-1}$ is strictly proper if and only if each column of $N(z)$ has degree less than the degree of the corresponding column of $D(z)$.*

Proof. Suppose first that $G(z) = N(z)D(z)^{-1}$ is strictly proper. We have

$$N(z) = G(z)D(z)$$

and if $N_i(z)$ and $D_i(z)$ denote the i -th columns of $N(z)$ and $D(z)$ respectively,

$$N_i(z) = G(z)D_i(z).$$

Since $G(z)$ is strictly proper, the degree of an element of $N_i(z)$ is less than the degree of the corresponding element of $D_i(z)$. (Note that for this part of the argument we did not need $D(z)$ to be column reduced.)

Assume now that $D(z)$ is column reduced and that the degree of each column of $N(z)$ is less than the degree of the corresponding column of $D(z)$. We may write

$$D(z) = HS(z) + L(z),$$

where H is the leading coefficient matrix of $D(z)$. Then

$$D(z)^{-1} = S(z)^{-1}H^{-1}(I + L(z)S(z)^{-1}H^{-1})^{-1}$$

Therefore

$$N(z)D(z)^{-1} = (N(z)S(z)^{-1})H^{-1}(I + L(z)S(z)^{-1}H^{-1})^{-1}$$

is the product of two rational matrices. The factor $N(z)S(z)^{-1}$ is strictly proper by hypothesis. Regarding the second factor, $L(z)S(z)^{-1}H^{-1}$ is strictly proper and so by Lemma 13.1.1, we see that $(I + L(z)S(z)^{-1}H^{-1})^{-1}$ is a proper rational matrix. It follows that $N(z)D(z)^{-1}$ is strictly proper, as required. \square

13.4.2 Theorem. Suppose $D(z)$ and $A(z)$ are polynomial matrices and $D(z)$ is invertible and column-reduced. Then there are unique polynomial matrices $Q(z)$ and $R(z)$ such that for each i , the degree of the i -th column of P_1 is less than the degree of the i -th column of D , and

$$A(z) = Q(z)D(z) + R(z).$$

Proof. The matrix $A(z)D(z)^{-1}$ is rational and so

$$A(z)D(z)^{-1} = Q(z) + P(z),$$

where $P(z)$ is polynomial and $R(z)$ is a strictly proper rational matrix. Hence

$$A(z) = Q(z)D(z) + P(z)D(z)$$

and, since $A(z)$ and $Q(z)D(z)$ are polynomial matrices, so is $P(z)D(z)$. Let $R(z) := P(z)D(z)$. Then $R(z)D(z)^{-1}$ is strictly proper and so by Theorem 13.4.1 the degree of each column of $R(z)$ has degree less than the degree of the corresponding column of $D(z)$.

Now suppose

$$A(z) = Q_1(z)D(z) + R_1(z)$$

where P_1 and Q_1 are polynomial and for each i , the degree of the i -th column of P_1 is less than the degree of the i -th column of D . Then

$$(Q - Q_1)D + (R - R_1) = 0$$

and therefore

$$Q - Q_1 = (R_1 - R)D^{-1}.$$

Here the left side is a polynomial matrix, while by Theorem 13.4.1, the right side is a strictly proper rational matrix. Therefore both sides are zero, and therefore $Q(z)$ and $R(z)$ are unique. \square

Note that we do not get a version of the Euclidean algorithm, because there is no guarantee that the remainder $R(z)$ is not a zero divisor or, if not, that it is reduced. So we cannot expect to be able to divide $Q(z)$ by $R(z)$.

13.5 Cayley-Hamilton

Suppose A is square. The matrix $zI - A$ is column reduced and linear so if we divide by it, the remainder must be a constant matrix. We can give an explicit formula for it.

13.5.1 Lemma. Suppose $F(z) = \sum_{i=0}^r F_i z^i$. Then remainder of $F(z)$ on right division by $zI - A$ is $\sum_i F_i A^i$.

Proof. We can write

$$(zI - A)^{-1} = z^{-1} \sum_{i \geq 0} A^i z^{-i};$$

the coefficient of z^{-1-j} in $F(z)(zI - A)^{-1}$ is then

$$F_0 A^j + F_1 A^{j+1} + \cdots + F_r A^{j+r} = (F_0 + F_1 A + \cdots + F_r A^r) A^j.$$

Therefore the strictly proper part of $F(z)(zI - A)^{-1}$ is

$$(zI - A)^{-1} (F_0 + A F_1 + \cdots + A^r F_r)$$

and the remainder on right division by $(zI - A)^{-1}$ is $F_0 + F_1 A + \cdots + F_r A^r$, as claimed. \square

This last result is an extension of the result that the remainder of $p(z)$ on division by $z - a$ is $p(a)$. It also implies the Cayley-Hamilton theorem. For suppose that $\phi(z)$ is the characteristic polynomial of A , and consider the remainder on left division of $\phi(z)I$ by $zI - A$. By the lemma, this remainder is $\phi(A)$. On the other hand

$$(zI - A)^{-1} \phi(z)I = \phi(z)^{-1} \text{adj}(zI - A) \phi(z) = \text{adj}(zI - A).$$

As $\text{adj}(zI - A)$ is a matrix polynomial, it follows that $\phi(A)$ is zero.

We write $F(A)$ to denote the remainder of $F(z)$ on right division by $zI - A$. It follows from our results that a polynomial $f(z)$ is satisfied by A if and only if there is a matrix polynomial $Q(z)$ such that

$$f(z)I = Q(z)(zI - A).$$

If $\phi(z)$ is the characteristic polynomial of A then

$$\phi(z)I = \text{adj}(zI - A)(zI - A), \quad (13.5.1)$$

which implies that $\phi(A) = 0$. This is the Cayley-Hamilton theorem.

Let $d(z)$ denote the greatest common divisor of the entries of $\text{adj}(zI - A)$ and let $C(z)$ be the matrix polynomial

$$C(z) = d(z)^{-1} \text{adj}(zI - A).$$

If $p(z)$ is the polynomial $\phi(z)/d(z)$ then, since

$$\phi(z)I = \text{adj}(zI - A)(zI - A)$$

we see that $p(A) = 0$. Let $\psi(z)$ be the minimal polynomial of A and let $\Psi(z)$ be the matrix polynomial satisfying

$$\psi(z)I = \Psi(z)(zI - A).$$

If $c(z) := p(z)/\psi(z)$, then

$$C(z)(zI - A) = p(z)I = c(z)\psi(z)I = c(z)\Psi(z)(zI - A).$$

As $zI - A$ is invertible, this implies that $C(z) = c(z)\Psi(z)$. Since the greatest common divisor of the entries of $C(z)$ is 1, it follows that $c(z) = 1$. Thus we have shown that $\psi(z) = \phi(z)/d(z)$.

13.6 Greatest Common Divisors

We say polynomial matrices $A(z)$ and $B(z)$ are *right equivalent* if there is a unimodular matrix $C(z)$ such that $A(z)C(z) = B(z)$.

If $A(z)$, $B(z)$ and $C(z)$ are polynomial matrices and $A(z)B(z) = C(z)$, we say that $A(z)$ is a *left divisor* and $B(z)$ a *right divisor* of $C(z)$. We do not insist that divisors be square, although this will be the most important case. We will only be interested in left divisors whose columns are linearly independent; equivalently those that have a rational left inverse. If the columns of A are linearly independent and $AX = B$ and $BY = A$, then $AXY = A$, and therefore XY is unimodular. Hence if A is a left divisor of B and B is a left divisor of A , then B and A are right equivalent.

We say that $D(z)$ is a *greatest common right divisor* of $A(z)$ and $B(z)$ if $D(z)$ is a right divisor of any right divisor of $A(z)$ and $B(z)$. If D and D_1 are two greatest common right divisors of A and B then

$$D_1 = XD, \quad D = YD_1$$

and therefore $D = YXD$. If the rows of D are linearly independent, it follows that $YX = I$ and therefore Y and X are unimodular. Consequently D and D_1 are left equivalent. Two polynomial matrices are *right coprime* if and only if all their greatest common right divisors are unimodular.

We show how to construct greatest common divisors. Let $D(z)$ and $N(z)$ be respectively $n \times n$ and $m \times n$ polynomial matrices. The Hermite normal form of

$$\begin{pmatrix} D(z) \\ N(z) \end{pmatrix}$$

is then

$$\begin{pmatrix} R(z) \\ 0 \end{pmatrix},$$

where R is $m \times m$ (and upper triangular). It follows that there is a unimodular 2×2 partitioned matrix such that

$$\begin{pmatrix} U_{1,1} & U_{2,1} \\ U_{2,1} & U_{2,2} \end{pmatrix} \begin{pmatrix} D \\ N \end{pmatrix} = \begin{pmatrix} R \\ 0 \end{pmatrix}.$$

It follows that

$$U_{1,1}D(z) + U_{2,1}N(z) = R(z). \quad (13.6.1)$$

On the other hand our 2×2 partitioned matrix is invertible, and so we have

$$\begin{pmatrix} D \\ N \end{pmatrix} = \begin{pmatrix} V_{1,1} & V_{2,1} \\ V_{2,1} & V_{2,2} \end{pmatrix} \begin{pmatrix} R \\ 0 \end{pmatrix}.$$

This implies that

$$D(z) = V_{1,1}R(z), \quad N(z) = V_{2,1}R(z),$$

and therefore $R(z)$ is a common right divisor of $D(z)$ and $N(z)$. If $S(z)$ is a common right divisor of $D(z)$ and $N(z)$, then it follows from (13.6.1) that $S(z)$ is a right divisor of $R(z)$. Therefore $R(z)$ is a greatest common divisor of $D(z)$ and $N(z)$.

Suppose A is an $n \times n$ matrix and B an $n \times m$ matrix over \mathbb{F} . The pair (A, B) is controllable if and only

$$\text{rk}(A - \lambda I \quad B) = n$$

for all complex numbers λ . It follows from the exercise*** below that this condition is satisfied if and only if A and B are left coprime. Similarly (C, A) is observable if and only if A and C are right coprime.

This section is one place where the module approach is very useful. Let us work with matrices over a ring R . If A and B are such matrices then $AX = B$ if and only if each column of B is an R -linear combination of the columns of

A. Thus A is a left divisor of B if and only if the R -module generated by the columns of B is contained in the R -module generated by the columns of A . If the columns of B and C lie in R^n , they are left coprime if and only if the columns of B and C together generate R^n .

(1) Prove that $D(z)$ and $N(z)$ are right coprime if and only if

$$\begin{pmatrix} D(z) \\ N(z) \end{pmatrix}$$

has a rational left inverse.

13.7 An Identity

We will need the following result.

13.7.1 Lemma. Let $C = (c_{i,j})$ be a square matrix. Then

$$\frac{\partial}{\partial c_{i,j}} (zI - C)^{-1} = (zI - C)^{-1} e_i e_j^T (zI - C)^{-1}.$$

Proof. This is an easy consequence of the following identity, which itself is easily verified.

$$(zI - C)^{-1} - (zI - D)^{-1} = (zI - C)^{-1} (C - D) (zI - D)^{-1}. \quad \square$$

The matrix Ψ in the proof of the next result is defined in ??.

13.7.2 Theorem. Let ψ be a polynomial of degree d , let C_ψ be its companion matrix and let $E_\psi(z)$ denote the $d \times d$ matrix with ij -entry equal to $\psi_i(z) z^{j-1} / \psi(z)$. Let N be the companion matrix of z^d . Then

$$(zI - C_\psi)^{-1} - E_\psi(z) = \left(N + zN^2 + \cdots + z^{d-2}N^{d-1} \right)^T.$$

Proof. The right side of this identity is independent of ψ (apart from its degree). By ??,

$$\psi(z)^{-1} \Psi(z) = (zI - C_\psi)^{-1} e_1$$

and therefore

$$E_\psi(z) = (zI - C_\psi)^{-1} e_1 \begin{pmatrix} 1 & z & \cdots & z^{d-1} \end{pmatrix}.$$

Our strategy is to show that $(zI - C_\psi)^{-1} - E_\psi(z)$ is independent of ψ , and then evaluate it when $\psi = z^d$.

Assume

$$\psi(z) = t^d + a_1 t^{d-1} + \cdots + a_d;$$

then $C_{i,d} = -a_i$ and, by the previous lemma

$$\frac{\partial}{\partial a_i} (zI - C_\psi)^{-1} = -(zI - C_\psi)^{-1} e_i e_d^T (zI - C_\psi)^{-1}.$$

From ?? we have

$$e_d^T (zI - C_\psi)^{-1} = \psi(z)^{-1} e_i (1 \quad z \quad \cdots \quad z^{d-1})$$

and therefore

$$\frac{\partial}{\partial a_i} (zI - C_\psi)^{-1} = -\psi(z)^{-1} (zI - C_\psi)^{-1} e_i (1 \quad z \quad \cdots \quad z^{d-1}).$$

We have

$$\frac{\partial}{\partial a_i} E_\psi(z) = \frac{\partial}{\partial a_i} (zI - C_\psi)^{-1} e_i (1 \quad z \quad \cdots \quad z^{d-1}).$$

Since

$$e_i (1 \quad z \quad \cdots \quad z^{d-1}) e_i (1 \quad z \quad \cdots \quad z^{d-1}) = e_i (1 \quad z \quad \cdots \quad z^{d-1}),$$

it follows that

$$\frac{\partial}{\partial a_i} (zI - C_\psi)^{-1} = \frac{\partial}{\partial a_i} E_\psi(z).$$

We conclude that $(zI - C_\psi)^{-1} - E_\psi(z)$ is independent of ψ .

Now suppose $\psi(z) = z^d$ and $N = C_\psi$. Then $N^d = 0$,

$$(zI - N)^{-1} = z^{-1} (I - z^{-1} N)^{-1} = \sum_{i=0}^{d-1} z^{-i} N^i,$$

and

$$(E_\psi(z))_{i,j} = z^{-i+j-1}.$$

The theorem follows at once. □

- (1) By setting $z = 0$ in Theorem 13.7.2, deduce the expression for the inverse of an invertible companion matrix in Theorem 3.4.1.

13.8 Resolvents

Let A be an $n \times n$ matrix. In this section we generally work over any field that contains all the eigenvalues of A . The *resolvent* $R(z)$ of A is the matrix $(zI - A)^{-1}$. As

$$R(z) = \frac{1}{\det(zI - A)} \operatorname{adj}(zI - A),$$

each entry of $R(z)$ is a rational function. Let θ be an eigenvalue of A with multiplicity m . Then there are matrices A_i such that

$$R(z) = \sum_{r=-m}^{\infty} A_r(z - \theta)^r;$$

we wish to determine these matrices.

The key to this is the following identity.

13.8.1 Theorem. *If $R(z)$ is the resolvent of some matrix then*

$$R(z) - R(w) = -(z - w)R(z)R(w).$$

Proof. Let $R(z)$ be the resolvent of A . Then

$$(zI - A)(R(z) - R(w))(wI - A) = (wI - A) - (zI - A) = (w - z)I,$$

whence the result follows. \square

We note a simple consequence of this.

13.8.2 Lemma. *If A is symmetric then all poles of the entries of $R(z)$ are simple.*

Proof. Suppose that θ is an eigenvalue of A and

$$R(z) = \sum_{r \geq -m} A_r(z - \theta)^r.$$

Here $m \geq 1$ and we may assume without loss that $A_{-m} \neq 0$. Then $R(z)^T R(z)$ is equal to $A_{-m}^T A_{-m}(z - \theta)^{2m}$, plus terms of higher order and, as A is symmetric, $R(z)^T R(z) = R(z)^2$. On the other hand, from Theorem 13.8.1 we have that

$$\frac{d}{dz} R(z) = -R(z)^2.$$

The term of least order in $R(z)'$ is $-mA_{-m}(z - \theta)^{m+1}$; consequently we must have $m + 1 = 2m$, i.e., $m = 1$. \square

13.8.3 Lemma. Suppose that $R(z)$ is the resolvent of A and that θ is an eigenvalue of A with multiplicity m . If $R(z) = \sum_{-m}^{\infty} A_r(z - \theta)^r$ then

$$A_r A_s = \begin{cases} -A_{r+s+1}, & r, s \geq 0; \\ A_{r+s+1}, & r, s \leq -1; \\ 0, & \text{otherwise.} \end{cases}$$

Proof. We assume that 0 is an eigenvalue of A , and seek to determine the coefficients A_r in the expansion $R(z) = \sum_{r \geq -m} A_r z^r$. From Theorem 13.8.1 we have

$$-\sum_{r, s \geq -m} A_r A_s z^r w^s = -R(z)R(w) = \frac{R(z) - R(w)}{z - w} = \sum_{r \geq -m} A_r \frac{z^r - w^r}{z - w}.$$

The lemma follows for $\theta = 0$ by comparing coefficients of $z^i w^j$ in the two series above, and the general result is an easy consequence of this. \square

From this result we see that the matrices A_i , $i = -m, -m + 1, \dots$ commute. We also find that:

$$\begin{aligned} A_r &= (-1)^r A_0^{r+1}, & \text{if } r \geq 0, \\ A_{-r} &= (A_{-2})^{r-1}, & \text{if } r \geq 2, \\ A_{-1} A_{-r} &= A_{-r}, & \text{if } r \geq 0. \end{aligned}$$

Therefore the coefficients in our Laurent series for $R(z)$ are determined by A_0 , A_{-1} and A_{-2} , where $(A_{-1})^2 = A_{-1}$ and $(A_{-2})^m = 0$. Thus A_{-1} is idempotent and A_{-2} is nilpotent, let us denote them respectively by E_θ and N_θ . Now note that

$$(tI - A)R(z) = ((t - z)I + zI - A)R(z) = (t - z)R(z) + I;$$

Putting $t = \theta$ in this yields

$$(\theta I - A)A_r = A_{r-1}, \quad r \neq 0, \quad (\theta I - A)A_0 = A_{-1} - I. \quad (\theta I - A)A_{-m} = 0.$$

Hence

$$N_\theta = (\theta I - A)E_\theta.$$

Define the *principal part* $P_\theta(z)$ of $R(z)$ by

$$P_\theta(z) := \sum_{r=1}^m A_{-r}(z - \theta)^{-r}.$$

Thus

$$\begin{aligned} P_\theta(z) &= (z - \theta)^{-1} E_\theta + \sum_{r=1}^{m-1} N_\theta^r (z - \theta)^{-r} \\ &= (z - \theta)^{-1} \sum_{r=0}^{m-1} (\theta I - A)^r E_\theta (z - \theta)^{-r} \end{aligned}$$

We note that, if θ and τ are distinct eigenvalues of A , then $P_\theta P_\tau = 0$ and so $E_\theta E_\tau = 0$. We have the following result, which provides a partial fraction decomposition of the resolvent.

13.8.4 Theorem. *Let $R(z)$ be the resolvent of A and let $P_\theta(z)$ be the principal part of $R(z)$ at θ . Then $R(z) = \sum_\theta P_\theta(z)$.*

Proof. A rational function in z is called *proper* if the degree of its numerator is less than the degree of its denominator. A proper rational function with no poles is constant. The set of proper rational functions is a vector space.

We note that the entries of $R(z)$ and the entries of $P_\theta(z)$ are proper rational functions. Hence each entry of the difference

$$R(z) - \sum_\theta P_\theta(z);$$

is a proper rational function. By the construction of $P_\theta(z)$, these rational functions have no poles. As both $R(z)$ and $P_\theta(z)$ converge to zero as $z \rightarrow \infty$, our theorem follows. \square

We know that, if m is the multiplicity of θ as an eigenvalue of A then $A_r = 0$ when $r < -m$, equivalently $(\theta I - A)^m E_\theta = 0$. This implies that the order of the pole of $R(z)$ at θ is at most m .

13.8.5 Theorem. *The order of the pole of $R(z)$ at θ is equal to the multiplicity of θ as a zero of the minimal polynomial of A .*

Proof. Let $\psi(z)$ denote the minimal polynomial of A , let $\nu(\theta)$ be the multiplicity of θ as a zero of $\psi(z)$ and suppose

$$\psi_\theta(z) = \frac{\psi(z)}{(z - \theta)^{\nu(\theta)}}.$$

Let \mathcal{A}_θ denote the space spanned by the matrices $(\theta I - A)^i E_\theta$ and let $d(\theta)$ be its dimension. Thus $d(\theta)$ is the greatest integer such that $(\theta I - A)^{d(\theta)-1} E_\theta \neq 0$.

As $(\theta I - A)^{d(\theta)} P_\theta(z) = 0$, it follows that

$$\prod_{\theta} (\theta I - A)^{d(\theta)} R(z) = 0.$$

Since $R(z)$ is invertible, this implies that

$$\prod_{\theta} (\theta I - A)^{d(\theta)} = 0.$$

From the definition of the minimal polynomial we then deduce that $v(\theta) \leq d(\theta)$, for all eigenvalues θ of A . We show next that $v(\theta) = d(\theta)$.

The matrices $(\theta I - A)^i E_\theta$ for $i = 0, 1, \dots, d(\theta) - 1$ form a basis for \mathcal{A}_θ . As $\psi_\theta(\theta) \neq 0$, it follows that the matrix representing the action of $\psi_\theta(A)$ relative to this basis is triangular, with non-zero diagonal entries. In particular, it is invertible. On the other hand, if $M \in \mathcal{A}_\theta$, then

$$0 = (\theta I - A)^{v(\theta)} \psi_\theta(A) M = \psi_\theta(A) (\theta I - A)^v v(\theta) M,$$

and this implies that $(\theta I - A)^v$ acts as the zero operator on \mathcal{A}_θ . It follows that $v(\theta) \geq d(\theta)$. \square

13.8.6 Corollary. *For each eigenvalue θ , the matrix E_θ is a polynomial in A .*

Proof. Since $zR(z) \rightarrow I$ and $zP_\theta(z) \rightarrow E_\theta$ as $z \rightarrow \infty$, Theorem 13.8.4 implies that

$$I = \sum_{\theta} E_\theta. \quad (13.8.1)$$

It follows from the proof of Theorem 13.8.5 that $\psi_\theta(A) E_\tau = 0$ if $\tau \neq \theta$, whence (13.8.1) yields that

$$(\theta I - A)^i \psi_\theta(A) = (\theta I - A)^i \psi_\theta(A) E_\theta.$$

Referring to the proof of Theorem 13.8.5 again, we see that $\psi_\theta(A) E_\theta$ lies in \mathcal{A}_θ . It is not hard to show that the matrices $(\theta I - A)^i \psi_\theta(A) E_\theta$ for $i = 0, 1, \dots, v(\theta)$ form a basis for \mathcal{A}_θ , and accordingly each matrix in \mathcal{A}_θ must be a polynomial in A . \square

13.8.7 Corollary. *Any square matrix A is the sum of a diagonalizable and a nilpotent matrix, each of which is a polynomial in A .*

Proof. As $E_\theta E_\tau = 0$ when $\theta \neq \tau$ and $E_\theta^2 = E_\theta$, the column space of E_θ is an eigenspace for all the idempotents E_τ . Given this, (13.8.1) implies that \mathbb{F}^n is the direct sum of eigenspaces of E_θ . Hence E_θ is diagonalizable; more generally any linear combination of the matrices E_θ is diagonalizable. It is also a polynomial in A .

As $AE_\theta = E_\theta + N_\theta$, it also follows from (13.8.1) that

$$A = \sum_{\theta} (\theta E_\theta + N_\theta) = \sum_{\theta} \theta E_\theta + \sum_{\theta} N_\theta.$$

Since $N_\theta N_\tau = 0$ when $\theta \neq \tau$, it follows that $\sum_{\theta} N_\theta$ is nilpotent. Since $N_\theta = (\theta I - A)E_\theta$, we see that N_θ is a polynomial in A and, therefore, $\sum_{\theta} N_\theta$ is too. \square

The last result implies that symmetric matrices are diagonalizable—if A is symmetric, so is any polynomial in A , but the only symmetric nilpotent matrix is the zero matrix. It is slightly more difficult to see that the only normal nilpotent matrix is the zero matrix; from this it follows that normal matrices are diagonalizable.

13.8.8 Corollary. *Let $\varphi(z)$ be the characteristic polynomial of A and let $g(z)$ be the greatest common divisor of the determinants of the $(n-1) \times (n-1)$ submatrices of $zI - A$. Then $\varphi(z)/g(z)$ is the minimal polynomial of A .*

Proof. Let θ be an eigenvalue of A , with multiplicity m , and let ν be its multiplicity as a zero of $\psi(z)$. Let $f_{i,j}(z)$ be the ij -minor $zI - A$. It follows from Theorem 13.8.5 that no entry of $R(z)$ has a pole of order greater than $\nu(\theta)$ at θ , and that some entry has a pole of this order at θ . In other words $(z - \theta)^{m-\nu}$ divides each polynomial $f_{i,j}(z)$, and divides one of these polynomials exactly. The result follows immediately. \square

Chapter 14

Determinants

The determinant is a function on square matrices which plays many roles. If A is a square matrix over \mathbb{R} , its determinant is a measure of ‘what A does to volume’. More precisely, if S is a region in \mathbb{R}^n with unit volume, then the volume of set of points

$$\{Ax : x \in S\}$$

is $|\det(A)|$. Because of this, the determinant plays an important role in integration of functions of several variables.

14.1 Permutations

Let Ω be a set. A *permutation* of Ω is a bijection from Ω to itself. The set of all permutations of Ω is called the *symmetric group* on Ω . If $|\Omega| = n$, then $|\text{Sym}(\Omega)| = n!$. We use $\text{Sym}(n)$ to denote the set of all permutations on some set of size n , usually $\{1, \dots, n\}$. If $i \in \Omega$ and $\sigma \in \text{Sym}(\Omega)$, then we denote the image of i under σ by i^σ .

Permutations of Ω are functions from Ω to Ω , so if ρ and σ are permutations, their *product* $\sigma\rho$ is defined by

$$i^{\sigma\rho} = (i^\sigma)^\rho.$$

This is again a permutation of Ω . As we will see, the order matters: usually $\sigma\rho \neq \rho\sigma$. Since a permutation is a bijection, it has an inverse. If $\sigma \in \text{Sym}(\Omega)$, we denote the inverse of σ by σ^{-1} . We have

$$\sigma\sigma^{-1} = \sigma^{-1}\sigma.$$

The identity mapping on Ω is a bijection; we call it the *identity permutation* and denote it by 1. Finally, if ρ , σ and τ are permutations of Ω , then

$$(\rho\sigma)\tau = \rho(\sigma\tau).$$

In other words, multiplication of permutations is associative.

If $\Omega = \{1, \dots, n\}$ and $\sigma \in \Omega$, we can specify σ by writing down the sequence

$$1^\sigma, 2^\sigma, \dots, n^\sigma.$$

This is sometimes called the *Cartesian form* of the permutation. There is a second useful way to present permutations, which we develop now. Suppose $i \in \Omega$ and consider the infinite sequence of elements

$$i, i^\sigma, i^{\sigma^2}, \dots$$

by successively applying σ . Since Ω is finite there are integers r and s such that $r < s$ and

$$i^{\sigma^r} = i^{\sigma^s}.$$

Then

$$i = i^{\sigma^s \sigma^{-r}} = i^{\sigma^{s-r}}.$$

This shows that $r = 0$ and that s is the least integer such that $i^{\sigma^s} = i$. Hence the elements

$$i, i^\sigma, \dots, i^{\sigma^{s-1}}$$

are distinct. We call the cyclic sequence

$$(i, i^\sigma, \dots, i^{\sigma^{s-1}})$$

the *cycle* of σ that contains i . We can view σ as rotating the elements of this cycle.

We consider an example. Suppose $n = 7$ and the Cartesian form of σ is

$$2\ 3\ 1\ 5\ 6\ 7\ 4.$$

Then the cycle of σ that contains 1 is

$$(123)$$

and the cycle of σ that contains 5 is

$$(5674).$$

We regard this as equal to each of the cycles

$$(4567), (6745), (7456).$$

The distinct cycles of Ω form a partition of Ω . Together they determine σ —we can specify σ by simply listing its cycles. In the example at hand we may write

$$\sigma = (123)(4567).$$

The order in which we list the cycles is irrelevant. This is the *cyclic form* of σ . A permutation may have cycles of length one; it is conventional to omit this from the cyclic form if the underlying set is clear. (The cyclic form of the identity permutation is often denoted by (1).) Note that i lies in a cycle of length one if and only if it is fixed by σ , that is, $i^\sigma = i$.

Each cycle of a permutation is a permutation in its own right, and a permutation is the product of the permutations corresponding to its cycles.

A permutation is a *transposition* if it has one cycle of length two, and all other cycles have length one.

14.1.1 Theorem. *If $\sigma \in \text{Sym}(n)$ and σ has exactly k cycles, then it is the product of $n - k$ transpositions.* \square

We leave the proof as an exercise. By way of a hint we note that

$$(1234) = (12)(13)(14),$$

from which we see that a cycle of length m is the product of $m - 1$ transpositions. We must count cycles of length one.

14.2 The Sign of a Permutation

A function of x_1, \dots, x_n is *alternating* if, when τ is a transposition in $\text{Sym}(n)$,

$$f^\tau = -f.$$

Thus $x_1 - x_2$ is an alternating function of two variables. If f is symmetric and g alternating in x_1, \dots, x_n , then fg is alternating. Define the function $V(x_1, \dots, x_n)$ by

$$V(x_1, \dots, x_n) = \prod_{i < j} (x_i - x_j).$$

Clearly V is alternating. Further, if $\sigma \in \text{Sym}(n)$, then

$$V^\sigma = \text{sign}(\sigma)V,$$

where $\text{sign}(\sigma) = \pm 1$. The value of $\text{sign}(\sigma)$ is called the *sign* of σ . If σ is a transposition, $\text{sign}(\sigma) = -1$.

14.2.1 Theorem. *If $\sigma, \tau \in \text{Sym}(n)$, then $\text{sign}(\sigma\tau) = \text{sign}(\sigma)\text{sign}(\tau)$.*

Proof. We have

$$V^{\sigma\tau} = (\text{sign}(\sigma)V)^\tau = \text{sign}(\sigma)\text{sign}(\tau)V$$

and therefore $\text{sign}(\sigma\tau) = \text{sign}(\sigma)\text{sign}(\tau)$. \square

By Theorem 14.1.1, each permutation is a product of transpositions, and therefore we have the following:

14.2.2 Corollary. *If f is an alternating function of n variables and $\sigma \in \text{Sym}(n)$, then $f^\sigma = \text{sign}(\sigma)f$.* \square

The set of even permutations is known as the *alternating group*.

Since each permutation is a product of cycles, if we know the sign of these cycles, we can use the previous lemma to get the sign of the permutation itself.

14.2.3 Lemma. *The sign of a cycle is odd if and only if its length is even.*

Proof. It follows from Theorem 14.1.1 that a cycle of length k can be written as the product of $k - 1$ transpositions. Since the sign of a transposition is odd, the sign of a cycle of length k is $(-1)^{k-1}$. \square

14.2.4 Corollary. *If a permutation has exactly e even cycles, its sign is $(-1)^e$.* \square

14.3 Permutation Matrices

Let \mathbb{F} be a field. If $\sigma \in \text{Sym}(n)$, let $P(\sigma)$ be the linear transformation that maps

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \mapsto \begin{pmatrix} x_{1^\sigma} \\ x_{2^\sigma} \\ \vdots \\ x_{n^\sigma} \end{pmatrix}.$$

Thus if e_1, \dots, e_n is the standard basis for $\mathbb{F}^{n \times 1}$, then $P(\sigma)$ maps e_j to $e_{j\sigma^{-1}}$. Hence the coordinate matrix for $P(\sigma)$ is

$$(e_{1\sigma^{-1}} \quad e_{2\sigma^{-1}} \quad \dots \quad e_{n\sigma^{-1}}).$$

The inverses are annoying, it may help to note that the i -th row of this matrix is $e_{i\sigma}^T$. We call $P(\sigma)$ a permutation operator and the matrix which represents it is a *permutation matrix*.

The product of two permutation operators is a permutation operator, and consequently the product of two permutation matrices is a permutation matrix.

If P is a permutation matrix then $PP^T = I$, and therefore $P^{-1} = P^T$.

A matrix is a permutation matrix if it is a 01-matrix, and exactly one entry in each row and column is equal to 1. We define a matrix to be a *monomial matrix* if there is at most one non-zero entry in each row and each column. It is not hard to verify that a matrix M is monomial if $M = PD$, where P is a permutation matrix and D is diagonal. Similarly DP is monomial. If P is a permutation matrix and D is diagonal, then

$$P^{-1}DP$$

is diagonal.

14.3.1 Lemma. *The product of two monomial matrices of the same order is a monomial matrix.*

Proof. Suppose P_1 and P_2 are permutation matrices and D_1 and D_2 are diagonal. Then P_1D_1 and P_2D_2 are monomial and

$$(P_1D_1)(P_2D_2) = P_1P_2(P_2^{-1}D_1P_2)D_2.$$

Here P_1P_2 is a permutation matrix and $(P_2^{-1}D_1P_2)D_2$ is a product of diagonal matrices, and so is diagonal. Hence $(P_1D_1)(P_2D_2)$ is a monomial matrix. \square

14.4 Definition of the Determinant

In this section we define the determinant of a square matrix, and develop some of its properties.

For this we will use a somewhat unusual matrix product: it is commutative and associative and distributes over addition. If A and B are $m \times n$ matrices, we define their *Schur product* $A \circ B$ by

$$(A \circ B)_{i,j} = A_{i,j} B_{i,j}.$$

There are no difficulties in working with this product. If A and P are $n \times n$ matrices and P is a permutation matrix, then $A \circ P$ is a monomial matrix.

The *determinant* is a function from the set of $n \times n$ matrices over a field (e.g., \mathbb{R} or \mathbb{C}) to the field itself. We define it in stages. If D is diagonal, then

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

If M is monomial, then $M = DP$ where D is diagonal and P is a permutation matrix. If $P = P(\sigma)$ for some permutation σ , we define $\text{sign}(P)$ to be $\text{sign}(\sigma)$ and then

$$\det(M) := \det(D) \text{sign}(P).$$

Note that

$$PD = (PDP^{-1})P$$

where PDP^{-1} is diagonal. Since PDP^{-1} is diagonal and $\det(PDP^{-1}) = \det(D)$,

$$\det(PD) = \det(PDP^{-1}) \text{sign}(P) = \det(D) \text{sign}(P).$$

It is implicit in this that, if P is a permutation matrix, then $\det(P) = \text{sign}(P)$.

To complete the definition of the determinant, let $\text{Perm}(n)$ denote the set of all $n \times n$ permutation matrices. If $A \in \text{Mat}_{n \times n}(\mathbb{F})$, we define

$$\det(A) := \sum_{\pi \in \text{Sym}(n)} \det(A \circ P(\pi)).$$

By way of example, if $n = 2$ then $\text{Perm}(2)$ consists of the two matrices

$$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

and so if

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix},$$

then

$$\det(A) = \det \begin{pmatrix} a & 0 \\ 0 & d \end{pmatrix} + \det \begin{pmatrix} 0 & b \\ c & 0 \end{pmatrix} = ad + (-1)bc = ad - bc.$$

14.4.1 Lemma. *Let A be an $n \times n$ matrix. If A is lower triangular, then*

$$\det(A) = \prod_{i=1}^n A_{i,i}.$$

Proof. Suppose $P \in \text{Perm}(n)$. If $\det(A \circ P) \neq 0$, then P must be lower triangular, but the identity matrix is the only lower triangular permutation matrix. Therefore $\det(A) = \det(A \circ I)$, and the lemma follows. \square

14.4.2 Lemma. *If A is a square matrix, $\det(A^T) = \det(A)$.*

Proof. We note first that if M is monomial, so is M^T . Further, if $M = DP$ where D is diagonal and P is a permutation matrix, then

$$\det(M^T) = \det(P^T D) = \det((P^T D P) P^{-1}),$$

Since $P^{-1} = P^T$, we see that $P^T D P$ is diagonal, and therefore

$$\det(M^T) = \det(P^T D P) \text{sign}(P^{-1}) = \det(D) \text{sign}(P) = \det(M).$$

Now

$$\begin{aligned} \det(A^T) &= \sum_{P \in \text{Perm}(n)} \det(A^T \circ P) \\ &= \sum_{P \in \text{Perm}(n)} \det(A \circ P^T)^T \\ &= \sum_{P \in \text{Perm}(n)} \det(A \circ P^T) \\ &= \sum_{P \in \text{Perm}(n)} \det(A \circ P) \\ &= \det(A). \end{aligned}$$

14.5 The Determinant is Multiplicative

The determinant is useful in particular because, if A and B are square matrices of the same order, then $\det(AB) = \det(A) \det(B)$. We work towards a proof of this.

We work with functions on $n \times n$ matrices. We may think of such a function δ as a function of n variables, the columns of the matrix. To indicate this, if A

is $n \times n$ and e_1, \dots, e_n is the standard basis of $\mathbb{F}^{n \times 1}$, we may use $\delta(Ae_1, \dots, Ae_n)$ in place of $\delta(A)$. A function $\delta : \text{Mat}_{n \times n}(\mathbb{F}) \rightarrow \mathbb{F}$ is *multilinear* if $\delta(A)$ is a linear function of each column of A . If δ is multilinear and $Ae_1 = x + y$, then

$$\delta(A) = \delta(Ae_1, \dots, Ae_n) = \delta(x, Ae_2, \dots, Ae_n) + \delta(y, Ae_2, \dots, Ae_n).$$

Note that trace, although it is a linear function of A , is not multilinear. However, if P is a permutation matrix then the function δ_P given by

$$\delta_P(A) = \det(A \circ P)$$

is multilinear. (Prove it.) If δ_1 and δ_2 are multilinear, then their sum, given by

$$(\delta_1 + \delta_2)(A) = \delta_1(A) + \delta_2(A),$$

is multilinear.

A function $\delta : \text{Mat}_{n \times n}(\mathbb{F}) \rightarrow \mathbb{F}$ is *alternating* if $\delta(A) = 0$ whenever two columns of A are equal. This usage is different from the one used in Section 14.2, but we will see that it is consistent with it.

We need two preliminary results.

14.5.1 Lemma. *If M_1 and M_2 are $n \times n$ monomial matrices, then $\det(M_1 M_2) = \det(M_1) \det(M_2)$.*

Proof. We may suppose that for $i = 1, 2$,

$$M_i = D_i P_i$$

where D_i is diagonal and P_i is a permutation matrix. Then

$$M_1 M_2 = D_1 P_1 D_2 P_2 = D_1 (P_1 D P_1^{-1}) P_1 P_2.$$

Here $P_1 D P_1^{-1}$ is diagonal, so $D_1 (P_1 D P_1^{-1})$ is diagonal and also $P_1 P_2$ is a permutation matrix. Therefore $M_1 M_2$ is monomial and

$$\begin{aligned} \det(M_1 M_2) &= \det(D_1 (P_1 D P_1^{-1})) \text{sign}(P_1 P_2) \\ &= \det(D_1) \det(D_2) \text{sign}(P_1) \text{sign}(P_2) \\ &= \det(D_1 P_1) \det(D_2 P_2) \\ &= \det(M_1) \det(M_2). \end{aligned}$$

This completes the proof. □

14.5.2 Lemma. *If A , B and P are $n \times n$ matrices and P is a permutation matrix, then $(A \circ B)P = (AP) \circ (BP)$.*

Proof. Suppose e_1, \dots, e_n is the standard basis and $Pe_i = e_j$. Then

$$\begin{aligned} ((AP) \circ (BP))e_i &= (AP)e_i \circ (BP)e_i \\ &= Ae_j \circ Be_j \\ &= (A \circ B)e_j \\ &= (A \circ B)Pe_i. \end{aligned}$$

Since this works for all i , we have proved the lemma. \square

14.5.3 Theorem. *The determinant is an alternating multilinear function of the columns of a matrix.*

Proof. Since the functions δ_P are multilinear and since \det is the sum of the functions δ_P , it follows that \det is multilinear.

To show that \det is alternating, we first prove that if Q is a permutation matrix, then $\det(AQ) = \det(A) \text{sign}(Q)$. Using the previous two lemmas, we have

$$\begin{aligned} \det(AQ) &= \sum_{P \in \text{Perm}(n)} \det((AQ) \circ P) \\ &= \sum_{P \in \text{Perm}(n)} \det[(A \circ PQ^{-1})Q] \\ &= \sum_{P \in \text{Perm}(n)} \det(A \circ (PQ^{-1})) \det(Q) \\ &= \det(Q) \sum_{P \in \text{Perm}(n)} \det(A \circ (PQ^{-1})). \end{aligned}$$

Since

$$\{P : P \in \text{Perm}(n)\} = \{PQ^{-1} : P \in \text{Perm}(n)\},$$

the last sum above equals $\det(A)$, we have proved that $\det(AQ) = \det(A) \text{sign}(Q)$, as claimed.

Now suppose columns i and j of A are equal, let τ be the transposition (ij) and let $T = P(\tau)$. Then $\text{sign}(T) = -1$, $T^2 = I$ and $AT = A$; hence

$$(A \circ P)T = (AT) \circ PT = A \circ PT$$

and consequently

$$\begin{aligned}\det(A \circ P) + \det(A \circ PT) &= \det(A \circ P) + \det((A \circ P)T) \\ &= \det(A \circ P) + \det(A \circ P) \det(T) \\ &= \det(A \circ P) - \det(A \circ P) \\ &= 0.\end{aligned}$$

The set $\{P, PT\}$ is the left coset of the subgroup $\{I, T\}$ of $\text{Perm}(n)$. For fixed T , the set $\text{Perm}(n)$ can be partitioned into pairs of the form $\{P, PT\}$ (prove this), and therefore it follows that $\det(A) = 0$. \square

One corollary of this proof is that if $P \in \text{Perm}(n)$, then $\det(AP) = \det(A) \text{sign}(P)$. Hence the determinant is an alternating function in the sense we used in Section 14.2. More generally, the same argument shows that if δ is an alternating function on $n \times n$ matrices and P is a permutation matrix, then

$$\delta(AP) = \delta(A) \text{sign}(P).$$

Therefore a function that is alternating in the sense of this section is alternating in the sense we used in Section 14.2, but the current definition is more useful if we work over fields such as \mathbb{Z}_2 .

Our next result is a converse to the previous theorem.

14.5.4 Theorem. *If δ is an alternating multilinear function on $n \times n$ matrices and $\delta(I) = 1$, then $\delta(A) = \det(A)$ for all $n \times n$ matrices.*

Proof. We have

$$Ae_j = \sum_{i=1}^n A_{i,j} e_i.$$

Since δ is multilinear,

$$\delta(A) = \delta(Ae_1, \dots, Ae_n) = \sum_{i=1}^n \delta(A_{i,1} e_i, Ae_2, \dots, Ae_n)$$

and, using even more subscripts,

$$\delta(A) = \sum_{1 \leq i_1, \dots, i_n \leq n} \delta(A_{i_1,1} e_{i_1}, \dots, A_{i_n,n} e_{i_n}). \quad (14.5.1)$$

Since δ is multilinear,

$$\delta(A_{i_1,1} e_{i_1}, \dots, A_{i_n,n} e_{i_n}) = \delta(e_{i_1}, \dots, e_{i_n}) \prod_{k=1}^n A_{i_k,k};$$

and since δ is alternating if $r < s$ and $i_r = i_s$, then

$$\delta(e_{i_1}, \dots, e_{i_n}) = 0$$

Hence in (14.5.1), the summands indexed by the sequences i_1, \dots, i_n that are not permutations are zero, and therefore

$$\delta(A) = \sum_{P \in \text{Perm}(n)} \delta(A \circ P).$$

This shows that δ is determined by the values it takes on monomial matrices.

If D is diagonal and P is a permutation matrix, then since δ is alternating,

$$\delta(DP) = \delta(D) \text{sign}(P).$$

Further, since δ is multilinear,

$$\delta(D) = \prod_{i=1}^n D_{i,i} \delta(I)$$

and therefore

$$\delta(DP) = \det(DP) \delta(I).$$

This completes the argument. □

14.5.5 Corollary. *If A and B are $n \times n$ matrices, then $\det(AB) = \det(A) \det(B)$.*

Proof. Consider the function δ from $\text{Mat}_{n \times n}(\mathbb{F})$ to \mathbb{F} , given by

$$\delta(B) := \det(AB).$$

It is easy to verify that this is alternating and multilinear, and therefore

$$\delta(B) = c_A \det(B)$$

for some scalar c_A . Taking $B = I$ in the definition of δ , we see that $c_A = \det(A)$ and therefore $\det(AB) = \det(A) \det(B)$. □

14.6 The Laplace Expansion

The determinant is remarkable for the number of different ways in which we can compute it. Here we describe an approach due to Laplace. You may be familiar with the case when $k = 1$, because this is the well-known expansion by cofactors.

If $T = \{t_1, \dots, t_k\}$, define $\|T\|$ by

$$\|T\| = \sum_{i=1}^k (t_i - i).$$

Let $A_{S,T}$ denote the submatrix of A with rows indexed by S and columns by T . If $|S| = |T| = 1$, then $A_{S,T}$ is just an entry of A . We use \bar{S} to denote the complement of S in $\{1, \dots, n\}$. Now we can state and prove a result known as Laplace's expansion of the determinant.

14.6.1 Theorem. *Let A be an $n \times n$ matrix and let S and S' be two subsets of $\{1, \dots, n\}$, with sizes k and $n - k$ respectively. Then*

$$\sum_{T:|T|=k} (-1)^{\|T\|} \det(A_{S,T}) \det(A_{S',\bar{T}}) = \begin{cases} (-1)^{\|S\|} \det(A), & \text{if } S' = \bar{S}; \\ 0, & \text{otherwise.} \end{cases}$$

Proof. We first consider the case where $S' = \bar{S}$. Let S and T be subsets of $\{1, \dots, n\}$ with size k . Then

$$\det A = \sum_T \sum_{\sigma: S^\sigma = T} \det(A \circ P(\sigma)).$$

Note that if σ maps S to T then it must map \bar{S} to \bar{T} . Hence

$$\sum_{\sigma: S^\sigma = T} \det(A \circ P(\sigma)) = (-1)^{\|T\|} \det(A_{S,T}) \det(A_{\bar{S},\bar{T}}).$$

Now suppose that $S' \neq \bar{S}$. Let A' be the matrix whose first k rows are the rows of A indexed by S_1 , and whose last $n - k$ rows are the rows of A indexed by S_2 . Since we know that Laplace's expansion holds when $S \cap S' = \emptyset$, we see that $\det(A')$ is equal to the sum on the left on the statement of the theorem. On the other hand, A' has a repeated row, and therefore $\det(A') = 0$. \square

Let $A(i|j)$ denote the matrix we get from the square matrix A by deleting row i and column j . Then $(-1)^{i+j} \det(A(i|j))$ is called the ij -cofactor of A . The following special case of the Laplace expansion is known the *expansion by cofactors* of $\det(A)$. This is sometimes used as a definition of the determinant.

14.6.2 Corollary. *Let A be an $n \times n$ matrix. Then*

$$\det(A) = (-1)^{i-1} \sum_{j=1}^n (-1)^{j-1} A_{i,j} \det(A(i|j)). \quad \square$$

Let A be an $n \times n$ matrix. We define the *adjugate* $\text{adj}(A)$ of A as follows:

$$\text{adj}(A)_{i,j} = (-1)^{i+j} \det A(i|j).$$

Thus if

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix},$$

then

$$\text{adj}(A) = \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}.$$

If

$$J = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$$

then $\text{adj}(J) = 0$.

Applying the previous theorem with $k = 1$, we obtain:

14.6.3 Corollary. *If A is a square matrix, then $A \text{adj}(A) = \det(A) I$.* □

It is also true that $\text{adj}(A)A = \det(A)I$; this can be proved using the transpose. We leave the proof as an exercise.

14.6.4 Corollary. *If A is a square matrix, then it is invertible if and only if $\det(A)$ is.*

Proof. If $\det(A)$ is invertible, the previous corollary implies that

$$A^{-1} = \det(A)^{-1} \text{adj}(A).$$

If A is invertible then

$$1 = \det(I) = \det(AA^{-1}) = \det(A) \det(A^{-1})$$

and therefore $\det(A)$ is invertible. □

The following identity is due to Jacobi.

14.6.5 Theorem. *Let A be an $n \times n$ matrix and suppose $S \subseteq \{1, \dots, n\}$. If $s = |S|$, then*

$$\det(\operatorname{adj}(A)_{\bar{S}, \bar{S}}) = \det(A)^{n-1-s} \det(A_{S,S}).$$

Proof. If M is $n \times n$, we have $\operatorname{adj}(M)M = \det(M)I$ and, taking determinants of both sides yields

$$\det(\operatorname{adj}(M)) \det(M) = \det(M)^n.$$

Therefore $\det(\operatorname{adj}(M)) = \det(M)^{n-1}$. Assume S consists of the first s elements of $\{1, \dots, n\}$. We have $\operatorname{adj}(A)A = \det(A)I$ whence $\operatorname{adj}(A)Ae_i = \det(A)e_i$ and

$$\operatorname{adj}(A) \begin{pmatrix} Ae_1 & \dots & Ae_s & e_{s+1} & \dots & e_n \end{pmatrix} = \begin{pmatrix} \det(A)I_s & ? \\ 0 & \operatorname{adj}(A)_{\bar{S}, \bar{S}} \end{pmatrix}$$

Taking the determinant of each side, we get

$$\det(A)^{n-1} \det(A_{S,S}) = \det(A)^s \det(\operatorname{adj}(A_{\bar{S}, \bar{S}})).$$

This yields the theorem. □

14.7 The Characteristic Polynomial of a Matrix

If A is a square matrix then $\det(tI - A)$ is a polynomial in t . It is called the *characteristic polynomial* of A . It is not too difficult to verify that if A is $n \times n$, then its characteristic polynomial is a monic polynomial of degree n . If

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

then

$$\det(tI - A) = t^2 - (a + c)t + (ac - bd).$$

The constant term of the characteristic polynomial of A is

$$\det(-A) = (-1)^n \det(A).$$

Suppose $A = LBL^{-1}$. Then

$$\begin{aligned} \det(tI - A) &= \det(tI - LBL^{-1}) = \det[L(tI - B)L^{-1}] \\ &= \det(B) \det(tI - B) \det(L^{-1}) \\ &= \det(tI - B). \end{aligned}$$

Thus we see that similar matrices have the same characteristic polynomial.

We leave the proof of the following as an exercise.

14.7.1 Lemma. *If $\phi(t)$ is the characteristic polynomial of the square matrix A , then the coefficient of t^{n-1} is $-\text{tr}(A)$.* \square

Our next result is called the Cayley-Hamilton theorem. Cayley proved it for 2×2 and 3×3 matrices.

14.7.2 Theorem. *If $\phi(t)$ is the characteristic polynomial of the square matrix A , then $\phi(A) = 0$.*

Proof. Each entry of $\text{adj}(tI - A)$ is a polynomial in t with degree at most $n - 1$. Hence there are matrices B_1, \dots, B_n such that

$$\text{adj}(tI - A) = B_n + tB_{n-1} + \cdots + t^{n-1}B_1$$

We want to show that each of the matrices B_1, \dots, B_n is a polynomial in A .

We have

$$\begin{aligned} (tI - A)\text{adj}(tI - A) &= t^n B_1 + t^{n-1}(B_2 - AB_1) + \cdots + t(B_n - AB_{n-1}) + (-A)B_n. \end{aligned} \quad (14.7.1)$$

Assume that

$$\phi(t) = t^n + a_1 t^{n-1} + \cdots + a_n.$$

From Corollary 14.6.3 we have

$$(tI - A)\text{adj}(tI - A) = (t^n + a_1 t^{n-1} + \cdots + a_n)I. \quad (14.7.2)$$

If we equate the coefficients of the powers of t , we obtain:

$$B_1 = I, \quad B_{i+1} = AB_i + a_i I \quad (i = 1, \dots, n-1)$$

whence

$$\begin{aligned} B_1 &= I \\ B_2 &= A + a_1 I \\ B_3 &= AB_2 + a_2 I = A^2 + a_1 A + a_2 I \end{aligned}$$

and, in general,

$$B_{k+1} = A^k + a_1 A^{k-1} + \cdots + a_k I.$$

Thus B_k is a polynomial of degree $k - 1$ in A .

From (14.7.1) and (14.7.2), we see that $a_n I = -AB_n$. So

$$\begin{aligned} 0 &= AB_n + a_n I = A(A^{n-1} + a_1 A^{n-2} + \cdots + a_{n-1} I) + a_n I \\ &= \phi(A). \end{aligned}$$

This completes the proof. □

It is tempting to argue that if we substitute A for t in the equation

$$(tI - A) \operatorname{adj}(tI - A) = \phi(t)I,$$

then $tI - A$ becomes zero, and therefore $\phi(A) = 0$. It is true that if $f(t)$ is a polynomial in t with coefficients in a field and $t - a$ divides $f(t)$, then $f(a) = 0$. It need not be true that if $f(t)$ and $f_1(t)$ are polynomials in t with matrices as coefficients and

$$(tI - A)f_1(t) = f(t)$$

then $f(A) = 0$. The basic problem is, for example, that if b is a scalar then

$$t^2 b = tbt = bt^2,$$

but if A and B are square matrices, then the products $A^2 B$, ABA and BA^2 can all be different.

14.8 An Algorithm

If we attempt to compute the determinant of a matrix in $\operatorname{Mat}_{n \times n}(\mathbb{Z})$ using our definition, we may be obliged to sum $n!$ products. This is already unpleasant when $n = 4$. There is a second algorithm using elementary row operations; the only disadvantage of this is that its intermediate stages often require the use of rational numbers, even though the final answer is an integer. (This is the algorithm usually taught.) We are going to describe a third algorithm that does not suffer from this disadvantage, and still runs in polynomial time.

Let A be an $m \times n$ matrix and suppose $k \leq m, n$. We construct an $(m+1-k) \times (n+1-k)$ matrix $D_k(A)$ from A as follows. If $k \leq r \leq m$ and $k \leq s \leq n$, there is a unique $k \times k$ submatrix of A that contains the rs -entry of A along with all entries

in the first $k - 1$ rows and columns. Define $D_k(A)_{r-k, s-k}$ to be the determinant of this submatrix. So $D_1(A) = A$ and if

$$A = \begin{pmatrix} a_{1,1} & a_{1,2} & a_{1,3} \\ a_{2,1} & a_{2,2} & a_{2,3} \\ a_{3,1} & a_{3,2} & a_{3,3} \end{pmatrix},$$

then

$$D_2(A) = \begin{pmatrix} a_{1,1}a_{2,2} - a_{1,2}a_{2,1} & a_{1,1}a_{2,3} - a_{2,1}a_{1,3} \\ a_{1,1}a_{3,2} - a_{1,2}a_{3,1} & a_{1,1}a_{3,3} - a_{3,1}a_{1,3} \end{pmatrix}.$$

If A is $n \times n$, then $D_n(A) = \det(A)$. For any matrix A , let $d_k(A)$ denote the determinant of the submatrix formed by the first k rows and columns; we assume $d_0(A) = 1$.

14.8.1 Lemma. *If A is an $m \times n$ matrix, then $D_2(D_k(A)) = d_{k-1}(A)D_{k+1}(A)$.*

Proof. We prove the result by induction on the size of A . Since $D_1(A) = A$, the lemma holds when $k = 1$ and we assume $k \geq 2$.

First we consider a special case. Suppose A is $(k+1) \times (k+1)$. Then $D_{k+1}(A) = \det(A)$ and

$$D_k(A) = \begin{pmatrix} \det(A(k+1|k+1)) & \det(A(k+1|k)) \\ \det(A(k|k+1)) & \det(A(k|k)) \end{pmatrix}$$

Therefore

$$\begin{aligned} \det(D_k(A)) &= \det(A(k|k)) \det(A(k+1|k+1)) \\ &\quad - \det(A(k+1|k)) \det(A(k|k+1)) \end{aligned}$$

and so if $S := 1, \dots, k - 1$, then

$$D_2(D_k(A)) = \det(\operatorname{adj}(A)_{\bar{S}, \bar{S}}).$$

By Jacobi's identity (Theorem 14.6.5),

$$\det(\operatorname{adj}(A)_{\bar{S}, \bar{S}}) = \det(A) \det(A_{S,S}) = d_{k-1} D_{k+1}(A).$$

Now we verify that the result follows from this special case. If $i \geq k$ we and B is the matrix we get by deleting the i -th row of A , then $D_k(B)$ is obtained from $D_k(A)$ by deleting its $(i + 1 - k)$ -th row. Since that $D_k(A^T) = D_k(A)$, a similar claim holds when we delete columns.

If $i, j \geq k+1$, then $(D_{k+1}(M))_{i-k, j-k}$ is the determinant of the submatrix M of A formed by the intersection of rows 1 through k and i with columns 1 through k and j . Since $d_{k-1}(A) = d_{k-1}(M)$, we have

$$\begin{aligned} d_{k-1}(A)D_{k+1}(A)_{i-k, j-k} &= d_{k-1}(M)D_{k+1}(M) \\ &= D_2(D_k(M)) \\ &= D_2(D_k(A))_{i-k, j-k} \end{aligned}$$

and so the result follows. □

The algorithm to compute $\det(A)$ runs as follows. The input is an $n \times n$ matrix A . We also use a scalar δ , which is initially set to 1.

1. If $n = 1$, then $\det(A) = A$; halt.
2. If the first row or column of A is zero, then $\det(A) = 0$; halt.
3. If necessary, swap two columns of A so that $A_{1,1} \neq 0$ and replace δ by $-\delta$.
4. Compute $\delta^{-1}D_2(A)$ and let $\delta = (A)_{1,1}$. Return to the first step with $\delta^{-1}D_2(A)$ in place of A .

After $n - 1$ steps of this kind, we obtain $D_n(A) = \det(A)$.

We give one example. If

$$A := \begin{pmatrix} x & -1 & 0 \\ -1 & x & -1 \\ 0 & -1 & x \end{pmatrix}$$

then

$$D_2(A) = \begin{pmatrix} x^2 - 1 & -x \\ -x & x^2 \end{pmatrix}$$

Since $d_1(A) = x$,

$$\det(A) = D_3(A) = x^{-1}(x^4 - 2x^2) = x^3 - 2x.$$

This algorithm is sometimes attributed to C. Dodgson, better known as Lewis Carroll.

14.9 Summary

The most useful facts are (c), (f) and (g). You are not required to know anything about the proofs of (f), (g), (h), (i) and (j). You might need to use them. Note that (d) and (e) together yield an algorithm for computing the determinant, since we can bring a matrix to triangular form by elementary row operations.

- (a) Permutations, sign of a permutation, permutation and monomial matrices.
- (b) Definition of determinant.
- (c) $\det(A^T) = \det(A)$
- (d) If A is triangular, $\det(A) = \prod_i A_{i,i}$.
- (e) Adding a scalar multiple of one row of A to another does not change $\det(A)$. Swapping rows changes the sign. Ditto for columns. If we get B from A by multiplying a column by c , then $\det(B) = c \det(A)$.
- (f) Multilinear and alternating functions on matrices, $\det(AB) = \det(A) \det(B)$.
- (g) The adjugate of a matrix, $A \operatorname{adj}(A) = \det(A)I$.
- (h) Cofactor expansion of $\det(A)$.
- (i) The Cayley-Hamilton theorem.
- (j) Bareiss algorithm.
- (k) When the products AB and BA are both defined, $\det(I - AB) = \det(I - BA)$.
- (l) Binet-Cauchy.
- (m) $\det(\exp(M)) = \exp(\operatorname{tr}(M))$.

(We did not treat the last three items.)

14.10 Groups

In this chapter we met the ‘symmetric group’ and the ‘alternating group’. As we continue with the course, we will meet other ‘groups’. For the sake of background information, we explain the terminology.

A *group* is a set G with a multiplication \circ defined on it. If $a, b \in G$, then $a \circ b$ denotes the product of a and b . (In many cases the elements of G are operations on some structure, and $a \circ b$ denotes “do a , then b ”.) The multiplication must satisfy the following axioms.

1. If $a, b \in G$, then $a \circ b \in G$.
2. If $a, b, c \in G$, then $(a \circ b) \circ c = a \circ (b \circ c)$.
3. There is an element θ in G such that $\theta \circ a = a$ for all a in G .
4. For each element $a \in G$, there is an element a^{-1} in G such that $a^{-1} \circ a = \theta$.

The first axiom states that G is closed under multiplication. The element θ is the *identity element* of the group. The element a^{-1} is the *inverse* of a . We do **not** assume that $a \circ b = b \circ a$; if this does hold for all a and b the group is *commutative* (or *abelian*).

One example of a group is the integers, with $+$ as the ‘multiplication’. A second example is the set of invertible $n \times n$ matrices over a field with the usual matrix multiplication.

We usually write ab in place of $a \circ b$ unless G is commutative, in which case we write $a + b$. We usually use 1 to denote the identity unless G is commutative, when we use 0 .

Suppose $a, x, y \in G$ and $ax = ay$. Then

$$x = 1x = (a^{-1}a)x = a^{-1}(ax) = a^{-1}(ay) = (a^{-1}a)y = 1y = y.$$

Thus in a group we may ‘cancel on the left’. Since

$$a^{-1}(a1) = (a^{-1}a)1 = 1^2 = 1 = a^{-1}a,$$

it follows (by left cancellation) that $a1 = a$ for all a . Since

$$(aa^{-1})a = a(a^{-1}a) = a1 = a = 1a$$

we also see that $aa^{-1} = 1$ for any a . Now if $xa = ya$, then

$$x = x1 = x(aa^{-1}) = (xa)a^{-1} = (ya)a^{-1} = y(aa^{-1}) = y1 = y;$$

therefore we may also cancel on the right.

A subset of G is a *subgroup* if it contains the inverse of each of its elements and is closed under multiplication. The alternating group is a subgroup of the symmetric group.

Finally we point out that a group is a set with three operations. A binary operation which, given (a, b) as input, returns $a \circ b$. A unary operation which, given a as input, returns a^{-1} . And a nullary operation which, given no input, returns the identity θ . (It may help to understand the last statement if you think of a button on a calculator labelled π —this takes no input and returns π .)

Chapter 15

Rings, Fields, Algebras

This chapter is meant to provide some background, to help you deal with linear algebra over fields other than \mathbb{Q} , \mathbb{R} and \mathbb{C} .

15.1 Rings

A ring R consists of a set R on which an addition operation $+$ is defined, such that $(R, +)$ is a commutative group; in addition there is an associative multiplication in R that satisfies the usual distributive laws relative to addition. The multiplication is usually denoted by juxtaposition, i.e., the product of a and b is denoted ab (and ab need not equal ba). We always assume that there is a multiplicative identity, denoted by 1 (so $1x = x1 = x$ for all x in R).

The canonical examples are \mathbb{Z} , \mathbb{Q} , \mathbb{R} , \mathbb{C} . Polynomials over \mathbb{Q} , \mathbb{R} or \mathbb{C} form a ring, and so do power series. Further, matrices with entries from a ring R form a ring which is not normally commutative. Continuous real functions on \mathbb{R} form a ring.

Rings were first introduced in number theory, but now it is somewhat unusual for a mathematician not to be working in the context of some ring.

As a general principle, any operation we can carry out on abelian groups can be carried out on rings. So we have subrings, products and homomorphisms/quotients. Somewhat surprisingly, subrings do not play a big role, except for ideals (which you can look up). Also finite rings seem to be less useful than finite groups.

15.2 Fields

A field is a ring in which every non-zero element has a multiplicative inverse. The canonical examples are \mathbb{Q} , \mathbb{R} , \mathbb{C} . We see that \mathbb{Z} is not a field, and the rings of polynomials we referred to above are not fields (although they can be used to construct fields). If \mathbb{F} is a field then $\mathbb{F}(t)$, the ring of rational functions with coefficients from \mathbb{F} is a field.

As just defined, the multiplication in a field need not be commutative. However all fields we need are commutative and so henceforth field means commutative field.

The integers modulo a prime p form a field \mathbb{Z}_p . We consider this in some detail. Strictly speaking, the elements of \mathbb{Z}_p are equivalence classes of integers, where integers m and n are equivalent, i.e., $m \equiv n$, if p divides $m - n$. Each equivalence class contains exactly one element from the set of integers

$$\{0, 1, \dots, p - 1\}$$

and so we can identify the equivalence classes with the members of this set. It is not too difficult to show that the equivalence classes form a ring, with addition mod p and multiplication mod p as its operations. In fact we can show that, for any positive integer n , the set \mathbb{Z}_n forms a ring. But if n is not a prime we can write $n = ab$ where a and b both greater than 1, and therefore $ab = 0$ in \mathbb{Z}_n . It follows that the equivalence class of a does not have a multiplicative inverse—if $xa = 1$ and $ab = 0$ then

$$0 = x(ab) = (xa)b = 1b = b.$$

Therefore if n is not prime, then \mathbb{Z}_n is not a field.

If p is a prime then each non-zero element of \mathbb{Z}_p does have a multiplicative inverse. For if $a \in \mathbb{Z}_p$ and $a \neq 0$, then the gcd of a and p is 1, and hence there are integers x and y such that

$$xa + yp = 1,$$

and therefore $xa = 1$. Thus we can find the multiplicative inverse of a using the Euclidean algorithm. We have been a little sloppy here: when we apply the Euclidean algorithm we are viewing a and p as integers, but we originally chose a to be a non-zero element of the ring \mathbb{Z}_p . To avoid this we should use some notation like $[a]$ to denote the equivalence class of a , but the sloppiness is easier, and traditional.

We can also construct fields from rings of polynomials. Let \mathbb{F} be a field and let $\mathbb{F}[t]$ denote the ring of polynomials with coefficients from \mathbb{F} . If $p(t)$ is a monic polynomial in $\mathbb{F}[t]$, define a relation \equiv on $\mathbb{F}[t]$ by declaring polynomials g and h to be equivalent if their difference is divisible by p . Then this is an equivalence relation and the equivalence classes form a ring. You may show that this ring is a field if and only if p is irreducible over \mathbb{F} (has no non-trivial factors).

If we take $\mathbb{F} = \mathbb{R}$ and $p(t) = t^2 + 1$, this construction produces a field isomorphic to the complex numbers. If $\mathbb{F} = \mathbb{Z}_2$ and $p(t) = t^2 + t + 1$, we obtain a field with four elements.

Exercise: Let \mathbb{E} be a field and let F be the subset of \mathbb{E} consisting of all the elements of \mathbb{E} we can get by adding 1 to itself any number of times. (By assumption, $0 \in F$; thus F is the additive subgroup of \mathbb{E} generated by 1.) Show that F is a ring. If $|F|$ is finite, prove that it is a prime, and deduce that it is a field.

15.3 Algebras

A ring R is an algebra over a field \mathbb{F} if R is a vector space over \mathbb{F} such that if $x, y \in R$ and $a \in \mathbb{F}$, then

$$(ax)y = x(ay) = a(xy).$$

If 1 is the multiplicative identity in R , then the set $\{a1 : a \in \mathbb{F}\}$ forms a subring of R that is isomorphic to \mathbb{F} . Each element of this subring commutes with each element of R (it lies in the *center* of R).

The term 'algebra' has changed its meaning over the years, and it still has more than one interpretation. As we have just defined it, every algebra contains a multiplicative unit, but, in analysis for example, this requirement can be dropped.

The set of $d \times d$ matrices over of field \mathbb{F} forms an algebra. More generally, the set of linear mappings of a vector space to itself is an algebra. The complex numbers are an algebra over the reals.

The *dimension* of an algebra is its dimension as a vector space over the underlying field.

Let M denote the subset of the algebra of 2×2 matrices over \mathbb{Q} consisting of the matrices of the form

$$\begin{pmatrix} a & 2b \\ b & a \end{pmatrix} = a \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} + b \begin{pmatrix} 0 & 2 \\ 1 & 0 \end{pmatrix}.$$

It is not hard to show that this set is a subspace of $\text{Mat}_{2 \times 2}(\mathbb{Q})$ and this it is closed under multiplication. Hence it is a *subalgebra* of $\text{Mat}_{2 \times 2}(\mathbb{Q})$, but you can also show that it is commutative and that every non-zero element is invertible. Therefore it is a field, isomorphic to the field usually denoted by $\mathbb{Q}(\sqrt{2})$.

If \mathbb{A} is an algebra of dimension d over \mathbb{F} and $M \in \mathbb{A}$, then the $d + 1$ powers I, M, \dots, M^d are linearly dependent, whence there is a polynomial f such that $f(M) = 0$. Consequently there is a monic polynomial ψ of least degree such that $\psi(M) = 0$. It is called the minimal polynomial of M and degree at most d .

Exercise: If \mathbb{A} is a finite-dimensional algebra over \mathbb{F} and $x \in \mathbb{A}$, show that multiplication by x is a linear mapping (over \mathbb{F}).

Exercise: If \mathbb{A} is a finite-dimensional algebra over \mathbb{F} , prove that \mathbb{A} is isomorphic to an algebra of matrices over \mathbb{F} .

Exercise: Suppose K, L, M are fields with $K \leq L \leq M$. Then L and M are algebras over K ; let ℓ and m respectively denote the dimensions of L and M over K . Prove that ℓ divides m .

Exercise: Let \mathbb{F} be a field. If S is a subspace of $\text{Mat}_{d \times d}(\mathbb{F})$ such that each non-zero element is invertible, prove that $\dim(S) \leq d$.

Exercise: If \mathbb{A} is a finite-dimensional algebra over a field \mathbb{F} and each non-zero element of \mathbb{A} is invertible, prove that the minimal polynomial of each non-zero element is irreducible over \mathbb{F} .